

# Improving Active Learning Performance Through the Use of Data Augmentation

Joao Fonseca, and Fernando Bacao

**Abstract**—Active Learning (AL) is a well-known technique to optimize data usage in training, through the interactive selection of unlabeled observations, out of a large pool of unlabeled data, to be labeled by a supervisor. Its focus is to find the unlabeled observations that, once labeled, will maximize the informativeness of the training dataset, therefore reducing data related costs. The literature describes several methods to improve the effectiveness of this process. Nonetheless, there is a paucity of research developed around the application of artificial data sources in AL. This paper proposes a new framework for AL, which relies on the effective use of artificial data. Our method uses a hyperparameter optimization component to improve the generation of artificial instances during the AL process, as well as an uncertainty-based data selection method for the data generation mechanism. We compare the proposed method to the standard framework along with another active learning method that uses data augmentation. The models' performance was tested using four different classifiers, two AL-specific performance metrics and three classification performance metrics over 15 different datasets. We demonstrate that the proposed framework, using data augmentation, significantly improves the performance of AL, both in terms of classification performance and data selection efficiency.

**Index Terms**—Active Learning, Data Augmentation, Oversampling

## 1 INTRODUCTION

THE importance of training robust ML models with minimal data requirements is substantially increasing [1, 2, 3]. Although the growing amount of valuable data sources and formats being developed and explored is affecting various domains [4], this data is often unlabeled. Only a tiny amount of the data being produced and stored can be helpful in supervised learning tasks. Additionally, it is essential to note that labeling data for specific Machine Learning (ML) projects is often difficult and expensive, especially when data-intensive ML techniques are involved (e.g., Deep Learning classifiers) [1]. In this scenario, labeling the full dataset becomes impractical, time-consuming and expensive. Two different ML techniques attempt to address this problem: Semi-Supervised Learning (SSL) and Active Learning (AL). Even though they address the same problem, the two follow different approaches. SSL focuses on observations with the most certain predictions, whereas AL focuses on observations with the least certain predictions [5].

SSL attempts to use a small, predefined set of labeled and unlabeled data to produce a classifier with superior performance. This method uses the unlabeled observations to help define the classifier's decision boundaries [6]. Simultaneously, the amount of labeled data required to reach a given performance threshold is also reduced. It is a particular case of ML because it falls between the supervised and unsupervised learning perspectives. AL, instead of optimizing the informativeness of an existing training set, expands the dataset to include the most informative and/or representative observations [7]. It is an iterative process where a supervised model is trained and simultaneously

identifies the most informative unlabeled observations to increase the performance of that classifier. The combination of SSL with AL has been explored in the past, achieving state-of-the-art results [8].

Several studies have pointed out the limitations of AL within an Imbalanced Learning context [9, 10]. With imbalanced data, AL approaches frequently have low performance, high computational time, or data annotation costs. Studies addressing this issue tend to adopt classifier-level modifications, such as the Weighted Extreme Learning Machine [9, 11, 12]. However, classifier or query function-level modifications (See Section 2.1) have limited applicability since a universally good AL strategy has not yet been found [7]. Other methods address imbalance learning by weighing the observations as the function of the observation's class imbalance ratio [13]. Alternatively, other techniques reduce the imbalanced learning bias by combining Informative and Representative-based query approaches (see Section 2.1) [14]. Another approach to deal with imbalanced data and data scarcity, in general, is generating synthetic data [15]. This approach has the advantage of being classifier-agnostic, it potentially reduces the imbalanced learning bias, and also works as a regularization method in data-scarce environments, such as AL implementations [16]. However, most recent studies improve the AL performance by modifying the design/choice of the classifier and query functions used.

The usage of data augmentation in AL is not new. The literature found on the topic (see Section 2.3) focuses on either image classification or Natural Language Processing and uses Deep Learning-based data augmentation to improve the performance of neural network architectures in AL. These methods, although showing promising results, represent a limited perspective of the potential of data augmentation in a real-world setting. First, using Deep Learning

• J. Fonseca and F. Bacao are with NOVA Information Management School, Universidade Nova de Lisboa, Lisbon, Portugal.  
E-mail: {jpfonseca, bacao}@novaims.unl.pt

in an iterative setting requires access to significant computational power. Second, these models tend to use sophisticated data augmentation methods, whose implementation may not be accessible to non-sophisticated users. Third, the studies found on the topic are specific to the domain, classifier, and data augmentation method. Consequently, the direct effect of data augmentation is unclear: these studies implement different neural network-based techniques for different classification problems, whose performance may be attributed to various elements within the AL framework.

In this study, we explore the effect of data augmentation in AL in a context-agnostic setting, along with two different data augmentation policies: oversampling (where the amount of data generated for each class equals the amount of data belonging to the majority class) and non-constant data augmentation policies (where the amount of data generated exceeds the amount of data belonging to the majority class in varying quantities) between iterations. We start by conceptualizing the AL framework and each of its elements, as well as the modifications involved to implement data augmentation in the AL iterative process. We argue that simple, non-domain specific data augmentation heuristics are sufficient to improve the performance of AL implementations, without the need to resort to deep learning-based data augmentation algorithms.

When compared to the standard AL framework, the proposed framework contains two additional components: the Generator and the Hyperparameter Optimizer. We implement a modified version of the Geometric Synthetic Minority Oversampling Technique (G-SMOTE) [17] as a data augmentation method with an optimized generation policy (explained in Section 2.2). We also propose a hyperparameter optimization module, which is used to find the best data augmentation policy at each iteration. We test the effectiveness of the proposed method in 15 datasets of different domains. We implement three AL frameworks (standard, oversampling and varying data augmentation) using four different classifiers, three different performance metrics and calculate two AL-specific performance metrics.

The remainder of this manuscript is structured as follows: Section 2 introduces relevant topics discussed in the paper and describes the related work. Section 3 elucidates the proposed method. Section 4 details the methodology of the study's experiment. Section 5 presents the results obtained from the experiment, as well as a discussion of these results. Section 6 presents the conclusions drawn from this study.

## 2 BACKGROUND

### 2.1 Active Learning

This paper focuses on pool-based AL methods as defined in [18]. The goal of AL models is to maximize the performance of a classifier,  $f_c$ , while annotating as least observations,  $x_i$ , as possible. They use a data pool,  $\mathcal{D}$ , where  $\mathcal{D} = \mathcal{D}_{lab} \cup \mathcal{D}_{pool}$  and  $|\mathcal{D}_{pool}| \gg |\mathcal{D}_{lab}|$ .  $\mathcal{D}_{pool}$  and  $\mathcal{D}_{lab}$  refer to the sets of unlabeled and labeled data, respectively.

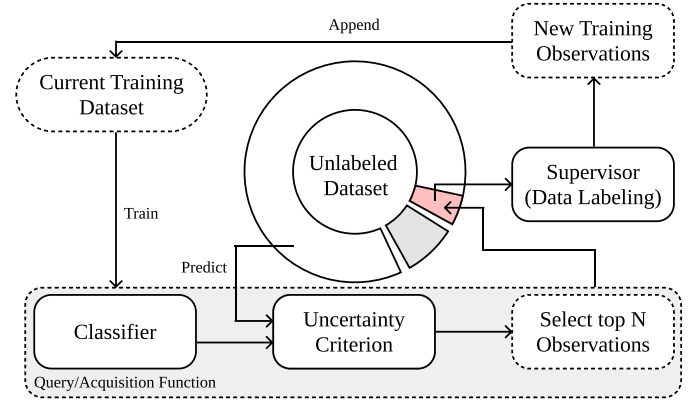


Fig. 1. Diagram depicting a typical AL iteration. In the first iteration, the training set collected during the initialization process becomes the "Current Training Dataset".

Having a budget of  $T$  iterations (where  $t = 1, 2, \dots, T$ ) and  $n$  annotations per iteration, at iteration  $t$ ,  $f_c$  is trained using  $\mathcal{D}_{lab}^t$  to produce, for each  $x_i \in \mathcal{D}_{pool}^t$ , an uncertainty score using an acquisition function  $f_{acq}(x_i; f_c)$ . These uncertainty scores are used to annotate the  $n$  observations with highest uncertainty from  $\mathcal{D}_{pool}^t$  to form  $\mathcal{D}_{new}^t$ . The iteration ends with the update of  $\mathcal{D}_{lab}^{t+1} = \mathcal{D}_{lab}^t \cup \mathcal{D}_{new}^t$  and  $\mathcal{D}_{pool}^{t+1} = \mathcal{D}_{pool}^t \setminus \mathcal{D}_{new}^t$  [19, 2]. This process is shown in Figure 1. Before the start of the iterative process, assuming  $\mathcal{D}_{lab}^{t=0} = \emptyset$ , the data used to populate  $\mathcal{D}_{lab}^{t=1}$  is typically collected randomly from  $\mathcal{D} = \mathcal{D}_{pool}^{t=0}$  and is labeled by a supervisor [20, 21, 22].

Research focused on AL has typically been focused on the specification of  $f_{acq}$  [23] and domain-specific applications. Acquisition functions can be divided into two different categories [24, 25]:

- 1) Informative-based. These strategies use the classifier's output to assess the importance of each observation towards the performance of the classifier [26].
- 2) Representative-based. These strategies estimate the optimal set of observations that will optimize the classifier's performance [25].

Although there are significant contributions toward the development of more robust query functions and classifiers in AL, modifications to AL's basic structure are rarely explored. In [21] the authors introduce a loss prediction module in the AL framework to replace the uncertainty criterion. This model implements a second classifier to predict the expected loss of the unlabeled observations (using the actual losses collected during the training of the original classifier) and return the unlabeled observations with the highest expected loss. However, this contribution is specific to neural networks (and more specifically, to deep neural networks) and was only tested for image classification.

### 2.2 Data Augmentation

Data Augmentation methods expand the training dataset by introducing new and informative observations [27]. The production of artificial data may be done via the introduction of perturbations on the input [28], feature [29],

or output space [27]. Data Augmentation methods may be divided into two categories [30]:

- 1) Heuristic approaches attempt to generate new and relevant observations by applying a predefined procedure, usually incorporating some degree of randomness [31]. Since these methods typically occur in the input space, they require fewer data and computational power when compared to Neural Network methods.
- 2) Neural Network approaches, on the other hand, map the original input space into a lower-dimensional representation, known as the feature space [29]. The generation of artificial data occurs in the feature space and is reconstructed into the input space. Although these methods allow the generation of less noisy data in high-dimensional contexts and more plausible artificial data, they are significantly more computationally intensive.

While some techniques may depend on the domain, others are domain-agnostic. For example, Random Erasing [32], Translation, Cropping and Flipping are examples of image data-specific augmentation methods. Other methods, such as autoencoders, may be considered domain agnostic.

### 2.3 Data Augmentation in Active Learning

The standard AL model can be complemented with a data augmentation function,  $f_{aug}(x_i; \tau)$ , where  $\tau$  defines the augmentation policy. In this context,  $\tau$  refers to the transformation applied and its hyperparameters and  $f_{aug}(x; \tau) : \mathcal{D} \rightarrow \mathcal{D}_{aug}(\mathcal{D})$  produces a modified observation,  $\tilde{x} \in \mathcal{D}_{aug}(\mathcal{D})$  where  $\mathcal{D}_{aug}(\mathcal{D})$  is the set of modified observations. This involves the usage of a new set of data,  $\mathcal{D}_{train}^t = \mathcal{D}_{lab}^t \cup \mathcal{D}_{aug}^t$ , to train the classifier.

As found in Section 2.1, improvements proposed in the AL framework are primarily focused on modifications of the classifier or query strategy. Furthermore, the few recent AL contributions implementing data augmentation were all (except one) applied to the computer vision or natural language processing (NLP) realm.

The only AL model found that uses data augmentation outside of the computer vision or NLP domains implements a pipelined approach, described in [20]. In this study, the AL model proposed is applied for tabular data using an oversampling data augmentation policy (*i.e.*, the artificial data was only generated to balance the target class frequencies). However, this AL model was applied in a Land Use/Land Cover context with specific characteristics that are not necessarily found in other supervised learning problems. Specifically, these types of datasets are high dimensional and have limited data variability within each class (*i.e.*, cohesive spectral signatures within classes) due to their geographical proximity. Furthermore, this method does not allow augmentation policy optimization (*i.e.*, every hyperparameter has to be hardcoded *a priori*).

The Bayesian Generative Active Deep Learning (BGDAL) [33] is another example of a pipelined combination of  $f_{acq}$  and  $f_{aug}$ , applied image classification. BGDAL uses a Variational AutoEncoder (VAE) architecture to generate artificial observations. However, the proposed model is computationally expensive, requires a large data pool to train the VAE, and is not only dependent on the quality of

the augmentations performed, but also on the performance of the discriminator and classifiers used.

The method proposed in [16], Look-Ahead Data Acquisition for Deep Active Learning, implements data augmentation to train a deep-learning classifier. However, adapting existing AL applications to use this approach is often impractical and implies the usage of image data since the augmentations used are image data specific and occur on the unlabeled observations, before the unlabeled data selection.

The Variational Adversarial Active Learning (VAAL) model [34] is a deep AL approach to image classification that uses as inputs the embeddings produced by a VAE into a secondary classifier, working as  $f_{acq}$ , to predict if  $x_i \in \mathcal{D}$  belongs to  $\mathcal{D}_{pool}$ . The  $n$  true positives with the highest certainty are labeled by the supervisor and  $\mathcal{D}_{pool}$  and  $\mathcal{D}_{lab}$  are updated as described in Section 2.1. The Task-aware VAAL model [35] extends the VAAL model by introducing a ranker, which consists of the Learning Loss module introduced in [21]. These models use data augmentation techniques to train the different neural network-based components of the proposed models. However, the AL components used are specific image classification, computationally expensive and the analysis of the effect of data augmentation in these AL models is not discussed.

In [36], the proposed AL method was explicitly designed for image data classification, where a deep learning model was implemented as a classifier, but its architecture is not described, the augmentation policies used are unknown and the results reported correspond to single runs of the discussed model. The remaining AL models found implement data augmentation for NLP applications, in [37, 38]. However, these methods were designed for specific applications within that domain and are not necessarily transferable to other domains or tasks.

## 3 PROPOSED METHOD

Based on the literature found on AL, most of the contributions and novel implementations of AL algorithms have focused on the improvement of the choice/architecture of the classifier or the improvement of the uncertainty criterion. In addition, the resulting classification performance of AL-trained classifiers is frequently inconsistent and marginally improve the classification performance when compared to classifiers trained over the entire training set. In addition, there is also significant variability in the data selection efficiency during different runs of the AL iterative process [20].

This paper provides a context-agnostic AL framework for the integration of Data Augmentation within AL, with the following contributions:

- 1) Improvement of the AL framework by introducing a parameter tuning stage only using the labeled dataset available at the current iteration (*i.e.*, no labeled hold-out set is needed).
- 2) Generalization of the generator module proposed in [20] from oversampling techniques to any other data augmentation mechanism and/or policy.
- 3) Implementation of data augmentation outside the Deep AL realm, which was not previously found in the literature.

- 4) Analysis of the impact of Data Augmentation and Over-sampling in AL over 15 different datasets of different domains, while comparing them with the standard AL framework.

The proposed AL framework is depicted in Figure 2. The generator element becomes an additional source of data and is expected to introduce additional data variability into the training dataset. This aspect should allow the classifier to generalize better and perform more consistently over unseen observations. However, in this scenario, the amount of data to generate per class at each iteration is unknown. Consequently, the hyperparameter tuning step was introduced to estimate the optimal data augmentation policy at each iteration. In our implementation, this step uses the current training dataset to perform an exhaustive search over specified generator parameters, tested over a 5-fold cross-validation method. The best augmentation policy found is used to train the iteration's classifier in the following step. This procedure is described in Algorithm 1.

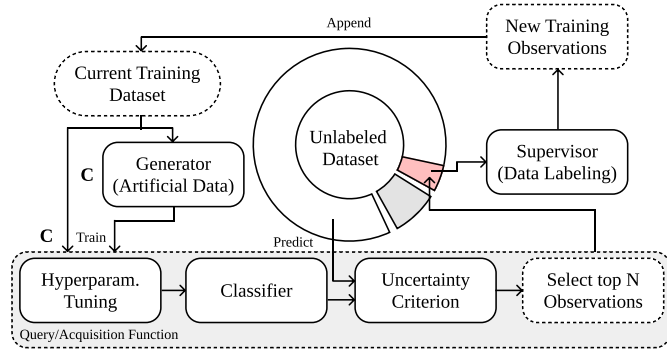


Fig. 2. Diagram depicting the proposed AL iteration. The proposed modifications are marked with a boldface “C.”

We implemented a simple modification in the selection mechanism of the G-SMOTE algorithm to show the effectiveness of data augmentation in an AL implementation. We use the uncertainties produced by  $f_{acq}$  to compute the probabilities of observations to be selected for augmentation as an additional parameter. This modification is described in Algorithm 2

This modification facilitates the usage of G-SMOTE beyond its original oversampling purposes. However, in this paper, the data augmentation strategies are also used to ensure that class frequencies are balanced. Furthermore, the amount of artificial data produced for each class is defined by the *augmentation factor*,  $\alpha_{af}$ , which represents a percentage of the majority class  $C_{maj}$  (e.g., an augmentation factor of 1.2 will ensure there are  $count(C_{maj}) \times 1.2$  observations in every class). In this paper's experiment, the data generation mechanism is similar to the one in [20]. This factor allows the direct comparison of the two frameworks and establishes a causality of the performance variations to the data generation mechanism (i.e., augmentation vs normal oversampling) and hyperparameter tuning steps. However, in this case, the hyperparameter tuning is solely going to be used for augmentation policy optimization.

In the proposed framework, we (1) generalize the generator module to accept any data augmentation method or

#### Algorithm 1: Proposed AL Framework (Single iteration)

---

**Given:**  $t \geq 1$ , performance metric  $f_{pm}$   
**Input:**  $\mathcal{D}_{pool}, \mathcal{D}_{lab}, f_c, f_{aug}, f_{acq}, \tau_{grid}, k, n$   
**Output:**  $\mathcal{D}_{pool}, \mathcal{D}_{lab}$

---

```

1 Function ParameterTuning( $f_c, f_{aug}, \tau_{grid}, \mathcal{D}_{lab}, k$ ):
2    $p \leftarrow 0$ 
3    $\tau \leftarrow \emptyset$ 
4    $\{\mathcal{D}_{lab}^1, \dots, \mathcal{D}_{lab}^k\} \leftarrow \mathcal{D}_{lab}$ 
      //  $\mathcal{D}_{lab}^n \cap \mathcal{D}_{lab}^m = \emptyset, \forall (n, m) \in 1, \dots, k$ 
5   forall  $\tau' \in \tau_{grid}$  do
6      $p' \leftarrow 0$ 
7     forall  $\mathcal{D}_{lab}^i \in \{\mathcal{D}_{lab}^1, \dots, \mathcal{D}_{lab}^k\}$  do
8        $\mathcal{D}'_{test} \leftarrow \mathcal{D}_{lab}^i$ 
9        $\mathcal{D}'_{train} \leftarrow \mathcal{D}_{lab} \setminus \mathcal{D}_{lab}^i$ 
10       $\mathcal{D}'_{train} \leftarrow f_{aug}(\mathcal{D}'_{train}; \tau')$ 
11      train  $f_c$  using  $\mathcal{D}'_{train}$ 
12       $p' \leftarrow p' \cup \{f_{pm}(f_c(\mathcal{D}'_{test}))\}$ 
13     $p' \leftarrow \frac{\sum_{x_i \in p'} x_i}{k}$ 
14    if  $p' > p$  then
15       $p \leftarrow p'$ 
16       $\tau \leftarrow \tau'$ 
17  return  $\tau$ 
18 begin
19   $\tau \leftarrow \text{ParameterTuning}(f_c, f_{aug}, \tau_{grid}, \mathcal{D}_{lab}, k)$ 
20   $\mathcal{D}_{train} \leftarrow f_{aug}(\mathcal{D}_{lab}; \tau)$ 
21  train  $f_c$  using  $\mathcal{D}_{train}$ 
22   $\mathcal{D}_{new} =$ 
       $\arg \max_{\mathcal{D}'_{pool} \subset \mathcal{D}_{pool}, |\mathcal{D}'_{pool}|=n} \sum_{x \in \mathcal{D}'_{pool}} f_{acq}(x; f_c)$ 
23  annotate  $\mathcal{D}_{new}$ 
24   $\mathcal{D}_{pool} \leftarrow \mathcal{D}_{pool} \setminus \mathcal{D}_{new}$ 
25   $\mathcal{D}_{lab} \leftarrow \mathcal{D}_{lab} \cup \mathcal{D}_{new}$ 

```

---

policy and (2) a hyperparameter tuning module to estimate the optimal data augmentation policy. This framework was designed to be task-agnostic. Specifically, any data augmentation method (domain-specific or not) may be applied, as well as any other parameter search method. It is also expected to be compatible with other AL modifications, including those that do not affect solely the classifier or uncertainty criterion, such as the one proposed in [21].

## 4 METHODOLOGY

This section describes the different elements included in the experimental procedure. The datasets used were acquired in open data repositories. Their sources and pre-processing steps are defined in Subsection 4.1. The classifiers used in the experiment are defined in Subsection 4.2. The metrics chosen to measure AL performance and overall classification performance are defined in Subsection 4.3. The experimental procedure is described in Subsection 4.4. The implementation of the experiment and resources used to do so are described in Subsection 4.5.

The methodology developed serves a two-fold purpose: (1) Compare classification performance once all the AL pro-

**Algorithm 2: G-SMOTE Modified for Data Augmentation in AL**


---

**Given:**  $t \geq 1, \mathcal{D}_{lab}^t \neq \emptyset, \mathcal{D}_{lab} = \mathcal{D}_{lab}^{min} \cup \mathcal{D}_{lab}^{maj}$ ,  
*GSMOTE*

**Input:**  $\mathcal{D}_{pool}^t, \mathcal{D}_{lab}^t, f_c^{t-1}, f_{acq}, \tau$

**Output:**  $\mathcal{D}_{train}^t$

```

1 Function DataSelection( $\mathcal{D}_{lab}^t, f_{acq}, f_c^{t-1}$ ):
2    $U \leftarrow \emptyset$ 
3    $P \leftarrow \emptyset$ 
4    $p_s \sim \mathcal{U}(0, 1)$ 
5   forall  $x_i \in \mathcal{D}_{lab}^t$  do
6      $u_{x_i} \leftarrow f_{acq}(x_i; f_c^{t-1})$ 
7      $U \leftarrow U \cup \{u_{x_i}\}$ 
8   forall  $u_{x_i} \in U$  do
9      $p_{x_i} \leftarrow \frac{u_{x_i}}{\sum U} + \sum P$ 
10     $P \leftarrow P \cup \{p_{x_i}\}$ 
11   $i \leftarrow \text{argmax}(P < p_s)$ 
12  return  $i$ -th element in  $\mathcal{D}_{lab}^t$ 
13 begin
14   $\mathcal{D}_{aug}^{min} \leftarrow \emptyset$ 
15   $\mathcal{D}_{aug}^{maj} \leftarrow \emptyset$ 
16   $\alpha_{af}, \alpha_{trunc}, \alpha_{def} \leftarrow \tau$ 
17   $N \leftarrow \text{count}(C_{maj}) \times \alpha_{af}$ 
18  forall  $\mathcal{D}'_{aug} \in \{\mathcal{D}_{aug}^{min}, \mathcal{D}_{aug}^{maj}\}$ ,
19     $\mathcal{D}'_{lab} \in \{\mathcal{D}_{lab}^{min}, \mathcal{D}_{lab}^{maj}\}$  do
20    while  $|\mathcal{D}'_{aug}| < N$  do
21       $x_{center} \leftarrow \text{DataSelection}(\mathcal{D}'_{lab}, f_{acq}, f_c^{t-1})$ 
22       $x_{gen} \leftarrow \text{GSMOTE}(x_{center}, \mathcal{D}'_{lab}, \alpha_{trunc}, \alpha_{def})$ 
23       $\mathcal{D}'_{aug} \leftarrow \mathcal{D}'_{aug} \cup \{x_{gen}\}$ 
24   $\mathcal{D}_{aug} \leftarrow \mathcal{D}_{aug}^{min} \cup \mathcal{D}_{aug}^{maj}$ 
25   $\mathcal{D}_{train}^t \leftarrow \mathcal{D}_{lab}^t \cup \mathcal{D}_{aug}$ 

```

---

cedures are completed (*i.e.*, optimal performance of a classifier trained via iterative data selection) and (2) Compare the amount of data required to reach specific performance thresholds (*i.e.*, the number of AL iterations required to reach similar classification performances).

#### 4.1 Datasets

The datasets used to test the proposed method are publicly available in open data repositories. Specifically, they were retrieved from OpenML and the UCI Machine Learning Repository. They were chosen considering diverse application domains, imbalance ratios, dimensionality and number of target classes, all of them focused on classification tasks. The goal is to demonstrate the performance of the different AL frameworks in various scenarios and domains. The data preprocessing approach was similar across all datasets. Table 1 describes the key properties of the 15 preprocessed datasets where the experimental procedure was applied.

TABLE 1

Description of the datasets collected after data preprocessing. The sampling strategy is similar across datasets. Legend: (IR) Imbalance Ratio

Dataset	Feat.	Inst.	Min. inst.	Maj. inst.	IR	Classes
Image Segmentation	14	1155	165	165	1.0	7
Mfeat Zernike	47	1994	198	200	1.01	10
Texture	40	1824	165	166	1.01	11
Waveform	40	1666	551	564	1.02	3
Pendigits	16	1832	176	191	1.09	10
Vehicle	18	846	199	218	1.1	4
Mice Protein	69	1073	105	150	1.43	8
Gas Drift	128	1987	234	430	1.84	6
Japanese Vowels	12	1992	156	323	2.07	9
Usps	256	1859	142	310	2.18	10
Gesture Segmentation	32	1974	200	590	2.95	5
Volkert	147	1943	45	427	9.49	10
Steel Plates	24	1941	55	673	12.24	7
Baseball	15	1320	57	1196	20.98	3
Wine Quality	11	1599	10	681	68.1	6

The data preprocessing pipeline is depicted as a flowchart in Figure 3. The missing values are removed from each dataset by removing the corresponding observations. This step ensures that the input data in the experiment is kept as close to its original form as possible. The non-metric features (*i.e.*, binary, categorical, and ordinal variables) were removed since the application of G-SMOTE is limited to continuous and discrete features. The datasets containing over 2000 observations were downsampled in order to maintain the datasets to a manageable size. The data sampling procedure preserves the relative class frequency of the dataset, in order to maintain the Imbalance Ratio (IR) originally found in each dataset (where  $IR = \frac{\text{count}(C_{maj})}{\text{count}(C_{min})}$ ). The remaining features of each dataset are scaled to the range of  $[-1, 1]$  to ensure a common range across features.

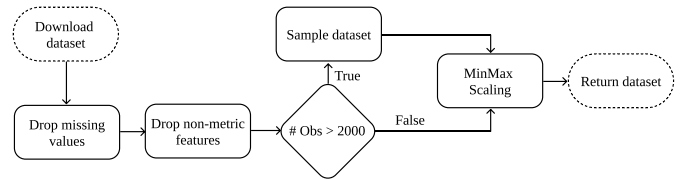


Fig. 3. Data preprocessing pipeline.

The preprocessed datasets were stored into an SQLite database file and is available along with the experiment's source code in the project's GitHub repository (see Subsection 4.5).

#### 4.2 Machine Learning Algorithms

We used a total of four classification algorithms and a heuristic data augmentation mechanism. The choice of classifiers was based on the popularity and family of the classifiers (tree-based, nearest neighbors-based, ensemble-based and linear models). Our proposed method was tested using a Decision Tree (DT) [39], a K-nearest neighbors classifier (KNN) [40], a Random Forest Classifier (RF) [41] and a Logistic Regression (LR) [42]. Since the target variables

are multi-class, the LR classifier was implemented using the one-versus-all approach. The predicted class is assigned to the label with the highest likelihood.

The oversampler G-SMOTE was used as a data augmentation method. The typical data generation policy of oversampling methods is to generate artificial observations on non-majority classes such that the number of majority class observations matches those of each non-majority class. We modified this data generation policy to generate observations for all classes, as a percentage of the number of observations in the majority class. In addition, the original G-SMOTE algorithm was modified to accept data selection probabilities based on classification uncertainty. These modifications are discussed in Section 3.

Every AL procedure was tested with different selection criteria: Random Selection, Entropy, and Breaking Ties. The baseline used is the standard AL procedure. As a benchmark, we add the AL procedure using G-SMOTE as a standard oversampling method, as proposed in [20]. Our proposed method was implemented using G-SMOTE as a data augmentation method to generate artificial observations for all classes, while still balancing the class distribution, as described in Section 3.

### 4.3 Evaluation Metrics

Considering the imbalanced nature of the datasets used in the experiment, commonly used performance metrics such as Overall Accuracy (OA), although being intuitive to interpret, are insufficient to quantify a model's classification performance [43]. The Cohen's Kappa performance metric, similar to OA, is also biased towards high-frequency classes since its definition is closely related to the OA metric, making its behavior consistent with OA [44]. However, these metrics remain popular choices for the evaluation of classification performance. Other performance metrics like  $Precision = \frac{TP}{TP+TN}$ ,  $Recall = \frac{TP}{TP+FN}$  or  $Specificity = \frac{TN}{TN+FP}$  are calculated as a function of True/False Positives (TP and FP) and True/False Negatives (TN and FN) and can be used on a per-class basis instead. In a multiple dataset scenario with varying amounts of target classes and meanings, comparing the performance of different models using these metrics becomes impractical.

Based on the recommendations found in [43, 45], we used two metrics found to be less sensitive to the class imbalance bias, along with OA as a reference for easier interpretability:

- The Geometric-mean scorer (G-mean) consists of the geometric mean of Specificity and Recall [45]. Both metrics are calculated in a multi-class context considering a one-versus-all approach. For multi-class problems, the G-mean scorer is calculated as its average per class values:

$$G\text{-mean} = \sqrt{Sensitivity \times Specificity}$$

- The F-score metric consists of the harmonic mean of Precision and Recall. The two metrics are also calculated considering a one-versus-all approach. The F-score for the multi-class case can be calculated using its average per class values [43]:

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

- The OA consists of the number of TP divided by the total amount of observations. Considering  $c$  as the label for the different classes present in a target class, OA is given by the following formula:

$$OA = \frac{\sum_c TP_c}{\sum_c (TP_c + FP_c)}$$

The comparison of the performance of AL frameworks is based on its data selection and augmentation efficacy. Specifically, an efficient data selection/generation policy allows the production of classifiers with high performance on unseen data while using as least non-artificial training data as possible. We follow the recommendations found in [46]. To measure the performance of the different AL setups, the performance of an AL setup will be compared using two AL-specific performance metrics:

- Area Under the Learning Curve (AULC). It is the sum of the classification performance over a validation/test set of the classifiers trained of all AL iterations. The resulting AULC scores are fixed within the range  $[0, 1]$  by dividing the AULC scores by the total amount of iterations (*i.e.*, the maximum performance area) to facilitate the interpretability of this metric.
- Data Utilization Rate (DUR) [47]. Measures the percentage of training data required to reach a given performance threshold, as a ratio of the percentage of training data required by the baseline framework. This metric is also presented as a percentage of the total amount of training data, without making it relative to the baseline framework. The DUR metric is measured at 45 different performance thresholds, ranging between  $[0.10, 1.00]$  at a 0.02 step.

### 4.4 Experimental Procedure

The evaluation of different active learners in a live setting is generally expensive, time-consuming, and prone to human error. Instead, a common practice is to compare them in an offline environment using labeled datasets [48]. Since the dataset is already labeled, the annotation process is done at zero cost in this scenario. Figure 4 depicts the experiment designed for one dataset over a single run.

A single run starts with the splitting of a preprocessed dataset into five different partitions, stratified according to the class frequencies of the target variable using the K-fold Cross Validation method. During this run, an active learner or classifier is trained five times using a different partition as the Test set each time. For each training process, a validation set containing 25% of the subset is created and is used to measure the data selection efficiency (*i.e.*, AULC and DUR using the classification performance metrics, specific to AL). Therefore, for a single training procedure, 20% of the original dataset is used as the validation set, 20% is used as the Test set and 60% is used as the training set. The AL simulations and the classifiers' training occur

within the training set. However, the classifiers used to find the maximum performance classification scores are trained over the full training set. The AL simulations are run over a maximum of 50 iterations (including the initialization step), adding 1.6% of the training set each time (*i.e.*, all AL simulations use less than 80% of the training set). Once the training phase is completed, the Test set classification scores are calculated using the trained classifiers. For the case of AL, the classifier with the optimal validation set score is used to estimate the AL's optimal classification performance over unseen data.

The process shown in Figure 4 is repeated over three runs using different random seeds over the 15 different datasets collected. The final scores of each AL configuration and classifier correspond to the average of the three runs and 5-fold Cross-Validation estimations (*i.e.*, the mean score of 15 fits, across 15 datasets).

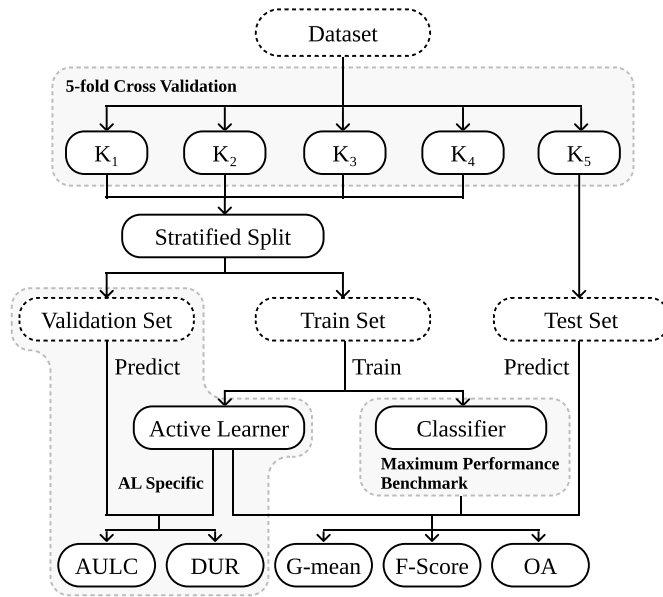


Fig. 4. Experimental procedure flowchart. The preprocessed datasets are split into five folds. One of the folds is used to test the best-found classifiers using AL and the classifiers trained using the entire training dataset (containing the remaining folds). The training set is used to run both the AL simulations as well as train the normal classifiers. The validation set is used to measure AL-specific performance metrics over each iteration. We use different AL-subsets for overall classification performance and AL-specific performance to avoid data leakage.

The hyperparameters defined for the AL frameworks, Classifiers, and Generators are shown in Table 2. In the Generators table, we distinguish the G-SMOTE algorithm working as a normal oversampling method from G-SMOTE-AUGM, which generates additional artificial data on top of the usual oversampling mechanism. Since the G-SMOTE-AUGM method is intended to be used with varying parameter values (via within-iteration parameter tuning), the parameters were defined as a list of various possible values.

#### 4.5 Software Implementation

The experiment was implemented using the Python programming language, along with the Python libraries Scikit-Learn [49], Imbalanced-Learn [50], Geometric-SMOTE [17],

Research-Learn and ML-Research libraries. All functions, algorithms, experiments, and results are provided on the project's GitHub repository.

## 5 RESULTS & DISCUSSION

In a multiple dataset experiment, the analysis of results should not rely upon the average performance scores across datasets uniquely. The domain of application and fluctuations of performance scores between datasets make the analysis of these averaged results less accurate. Instead, it is generally recommended to use the mean ranking scores to extend the analysis [51]. Since mean performance scores are still intuitive to interpret; we will present and discuss both results. The rank values are assigned based on the mean scores of three different 5-fold Cross-Validation runs (15 performance estimations per dataset) for each combination of dataset, AL configuration, classifier, and performance metric.

### 5.1 Results

The average rankings of the AL methods' AULC estimations are shown in Table 3. The proposed method almost always improves AL performance and ensures higher data selection efficiency.

Table 4 shows the average AULC scores, grouped by the classifier, Evaluation Metric and AL framework. The performance of the proposed method is almost always superior when considering the F-score and G-mean. On some occasions, the average AULC score is significantly improved when compared with the oversampling AL method.

The average DUR scores were calculated for various G-mean thresholds, varying between 0.1 and 1.0 at a 0.02 step (45 different thresholds in total). Table 5 shows the results obtained for these scores starting from a G-mean score of 0.6 and was filtered to show the thresholds ending with 0 or 6 only. In most cases, the proposed method reduces the amount of data annotation required to reach each G-mean score threshold.

The DUR scores relative to the Standard AL method are shown in Figure 5. A DUR below 1 means that the Proposed/Oversampling method requires less data than the Standard AL method to reach the same performance threshold. For example, running an AL simulation using the KNN classifier requires 80.7% of the amount of data required by the Standard AL method using the same classifier to reach an F-Score of 0.62 (*i.e.*, requires 19.3% less data).

The comparison of mean optimal classification scores of AL methods with Classifiers (using the entire training set, without AL) is shown in Table 6. Aside from the case of overall accuracy, the proposed AL method produces classifiers that almost consistently outperform classifiers using the whole training set (*i.e.*, the ones labeled as MP).

### 5.2 Statistical Analysis

When checking for statistical significance in a multiple dataset context it is critical to account for the multiple



TABLE 2  
Hyperparameter definition for the active learners, classifiers, and generators used in the experiment.

Active Learners	Hyperparameters	Inputs
Standard	# initial obs. # additional obs. per iteration max. iterations + initialization evaluation metrics selection strategy within-iteration param. tuning generator	1.6% 1.6% 50 G-mean, F-score, OA Random, Entropy, Breaking Ties None None
Oversampling	generator	G-SMOTE
Proposed	generator within-iteration param. tuning	G-SMOTE-AUGM Grid Search K-fold CV
Classifier		
DT	min. samples split criterion	2 gini
LR	maximum iterations multi-class solver penalty	100 One-vs-All liblinear L2 (Ridge)
KNN	# neighbors weights metric	5 uniform euclidean
RF	min. samples split # estimators criterion	2 100 gini
Generator		
G-SMOTE	# neighbors deformation factor truncation factor	4 0.5 0.5
G-SMOTE-AUGM	# neighbors deformation factor truncation factor augmentation factor	3, 4, 5 0.5 0.5 [1.1, 2.0] at 0.1 step

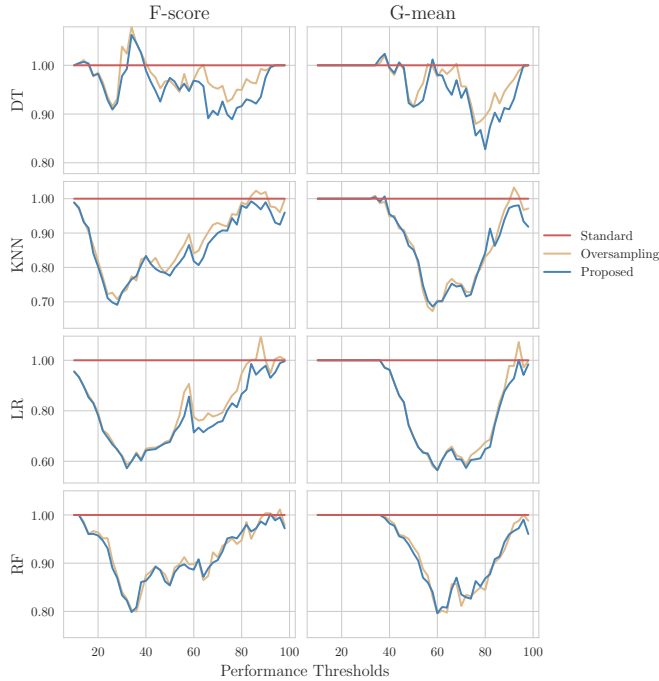


Fig. 5. Mean data utilization rates. The y-axis shows the percentage of data (relative to the baseline AL framework) required to reach the different performance thresholds.

TABLE 3  
Mean rankings of the AULC metric over the different datasets (15), folds (5), and runs (3) used in the experiment. The proposed method constantly improves the results of the original framework and, on average, almost always improves the results of the oversampling framework.

Classifier	Metric	Standard	Oversampling	Proposed
DT	Accuracy	2.13 ± 0.96	2.40 ± 0.49	<b>1.47 ± 0.62</b>
DT	F-score	2.47 ± 0.81	2.20 ± 0.40	<b>1.33 ± 0.70</b>
DT	G-mean	2.73 ± 0.57	1.93 ± 0.44	<b>1.33 ± 0.70</b>
KNN	Accuracy	2.07 ± 0.93	2.07 ± 0.68	<b>1.87 ± 0.81</b>
KNN	F-score	2.47 ± 0.81	1.87 ± 0.50	<b>1.67 ± 0.87</b>
KNN	G-mean	2.87 ± 0.34	<b>1.47 ± 0.50</b>	1.67 ± 0.70
LR	Accuracy	2.13 ± 0.88	2.20 ± 0.65	<b>1.67 ± 0.79</b>
LR	F-score	2.80 ± 0.40	1.87 ± 0.50	<b>1.33 ± 0.70</b>
LR	G-mean	2.80 ± 0.40	1.80 ± 0.54	<b>1.40 ± 0.71</b>
RF	Accuracy	2.27 ± 0.85	<b>1.87 ± 0.50</b>	<b>1.87 ± 0.96</b>
RF	F-score	2.73 ± 0.57	1.80 ± 0.54	<b>1.47 ± 0.72</b>
RF	G-mean	2.87 ± 0.34	<b>1.53 ± 0.50</b>	1.60 ± 0.71

comparison problem. Consequently, our statistical analysis focuses on the recommendations found in [51]. Overall, we perform three statistical tests. The Friedman test [52] is used to understand whether there is a statistically significant difference in performance between the three AL frameworks. As *post hoc* analysis, the Wilcoxon signed-rank test [53] was utilized to check for statistical significance between the performance of the proposed AL method and the oversampling AL method across datasets. As a second *post hoc*



TABLE 4

Average AULC of each AL configuration tested. Each AULC score is calculated using the performance scores of each iteration in the validation set. By the end of the iterative process, each AL configuration used a maximum of 80% instances of the 60% instances that compose the training sets (*i.e.*, 48% of the entire preprocessed dataset).

Classifier	Metric	Standard	Oversampling	Proposed
DT	Accuracy	0.663 ± 0.149	0.658 ± 0.153	<b>0.664 ± 0.155</b>
DT	F-score	0.610 ± 0.176	0.612 ± 0.179	<b>0.618 ± 0.181</b>
DT	G-mean	0.744 ± 0.129	0.751 ± 0.127	<b>0.755 ± 0.129</b>
KNN	Accuracy	<b>0.741 ± 0.160</b>	0.730 ± 0.178	0.734 ± 0.179
KNN	F-score	0.678 ± 0.208	0.684 ± 0.211	<b>0.687 ± 0.213</b>
KNN	G-mean	0.786 ± 0.152	<b>0.804 ± 0.139</b>	<b>0.804 ± 0.141</b>
LR	Accuracy	<b>0.736 ± 0.152</b>	0.723 ± 0.185	0.731 ± 0.184
LR	F-score	0.644 ± 0.228	0.673 ± 0.220	<b>0.682 ± 0.221</b>
LR	G-mean	0.767 ± 0.162	0.811 ± 0.134	<b>0.814 ± 0.136</b>
RF	Accuracy	<b>0.789 ± 0.148</b>	0.786 ± 0.153	0.785 ± 0.156
RF	F-score	0.724 ± 0.214	<b>0.735 ± 0.204</b>	<b>0.735 ± 0.205</b>
RF	G-mean	0.818 ± 0.150	<b>0.834 ± 0.135</b>	0.833 ± 0.135

TABLE 5

AL algorithms' mean data utilization as a percentage of the training set.

G-mean	Classifier	Standard	Oversampling	Proposed
0.60	DT	19.8%	<b>18.9%</b>	19.3%
0.60	KNN	18.4%	<b>11.8%</b>	12.8%
0.60	LR	23.0%	<b>9.7%</b>	<b>9.7%</b>
0.60	RF	14.1%	<b>7.7%</b>	7.8%
0.66	DT	23.1%	23.3%	<b>22.9%</b>
0.66	KNN	23.9%	<b>21.7%</b>	21.9%
0.66	LR	25.6%	<b>20.5%</b>	<b>20.5%</b>
0.66	RF	22.0%	17.6%	<b>17.5%</b>
0.70	DT	25.5%	25.0%	<b>24.8%</b>
0.70	KNN	26.8%	24.1%	<b>23.9%</b>
0.70	LR	29.9%	23.6%	<b>23.4%</b>
0.70	RF	23.8%	<b>22.1%</b>	22.3%
0.76	DT	33.4%	30.5%	<b>30.1%</b>
0.76	KNN	34.0%	27.7%	<b>27.3%</b>
0.76	LR	38.0%	27.6%	<b>26.2%</b>
0.76	RF	28.2%	<b>24.5%</b>	24.7%
0.80	DT	48.2%	43.8%	<b>41.2%</b>
0.80	KNN	38.8%	<b>34.4%</b>	34.6%
0.80	LR	43.7%	32.6%	<b>31.3%</b>
0.80	RF	32.4%	<b>27.2%</b>	27.7%
0.86	DT	69.6%	66.5%	<b>64.8%</b>
0.86	KNN	53.9%	<b>52.0%</b>	52.5%
0.86	LR	48.7%	45.3%	<b>45.0%</b>
0.86	RF	43.9%	<b>40.0%</b>	<b>40.0%</b>
0.90	DT	81.2%	79.4%	<b>76.6%</b>
0.90	KNN	60.9%	61.1%	<b>60.4%</b>
0.90	LR	62.1%	62.9%	<b>59.9%</b>
0.90	RF	57.1%	<b>55.7%</b>	56.2%
0.96	DT	100.0%	<b>99.7%</b>	100.0%
0.96	KNN	82.4%	79.7%	<b>77.1%</b>
0.96	LR	86.5%	84.0%	<b>81.8%</b>
0.96	RF	70.8%	71.1%	<b>70.3%</b>

analysis, the Holm-Bonferroni [54] method was employed to check for statistical significance between the methods using data generators and the Standard AL framework across classifiers and evaluation metrics.

Table 7 displays the *p-values* obtained with the Friedman test. The difference in performance across AL frameworks is statistically significant at a level of  $\alpha = 0.05$  regardless of the classifier or evaluation metric being considered.

Table 8 contains the *p-values* obtained with the Wilcoxon signed-rank test. The proposed method was able to outperform both the standard AL framework, as well as the AL

framework using a typical oversampling policy with statistical significance in 14 and 12 out of 15 datasets, respectively.

The *p-values* shown in Table 9 refer to the results of the Holm-Bonferroni test. The proposed method's superior performance was statistically significant in 9 out of 12 cases.

### 5.3 Discussion

In this paper, we study the application of data augmentation methods through the modification of the standard AL framework. This is done to further reduce the amount of labeled data required to produce a reliable classifier, at the expense of artificial data generation.

In Table 3, we found that the proposed method was able to outperform the Standard AL framework in all scenarios. Except for the overall accuracy metric, the mean rankings are consistent with the mean AULC scores found in Table 4, while showing performance improvements between the proposed method and both the standard and oversampling methods. The Friedman test in Table 7 showed that the difference in the performance of these AL frameworks are statistically significant, regardless of the classifier or performance metric being used.

The proposed method evidenced more consistent data utilization requirements in most of the assessed G-mean score thresholds when compared to the remaining AL methods, as seen in Table 5. For example, to reach a G-mean score of 0.9 using the KNN and LR classifiers, the average amount of data required with the Oversampling AL approach increased when compared to the standard approach. However, the proposed method was able to decrease the amount of data required in both situations. The robustness of the proposed method is clearer in Figure 5. In most cases, this method was able to outperform the Oversampling method. At the same time, the proposed method also addresses inconsistencies in situations where the Oversampling method was unable to outperform the standard method.

The statistical analyses found in Tables 8 and 9 revealed that the proposed method's superiority was statistically significant in all datasets except three (Baseball, Usps, and Volkert) and established statistical significance when compared to the standard AL method for all combinations of classifier and performance metric, except for three cases regarding the use of the overall accuracy metric. These results show that the proposed method increased the reliability of the new AL framework and improved the quality of the final classifier while using fewer data.

Even though it was not the core purpose of this study, we found that the proposed AL method consistently outperformed the maximum performance threshold. Specifically, in Table 6, the performance of the classifiers originating from the proposed method was able to outperform classifiers trained using the full training dataset in 9 out of 12 scenarios. This outcome suggests that the selection of a meaningful training subset training dataset paired with data augmentation not only matches the classification performance of ML algorithms, as it also improves them. Even in a setting with fully labeled training data, the proposed method may be used as a preprocessing technique to further optimize classification performance.

TABLE 6

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

Classifier	Metric	MP	Standard	Oversampling	Proposed
DT	Accuracy	<b>0.732 ± 0.155</b>	0.726 ± 0.157	0.721 ± 0.167	0.727 ± 0.168
DT	F-score	0.682 ± 0.194	0.679 ± 0.193	0.679 ± 0.197	<b>0.684 ± 0.200</b>
DT	G-mean	0.792 ± 0.138	0.791 ± 0.136	0.797 ± 0.134	<b>0.800 ± 0.137</b>
KNN	Accuracy	<b>0.801 ± 0.164</b>	0.799 ± 0.168	0.784 ± 0.183	0.789 ± 0.183
KNN	F-score	0.742 ± 0.224	0.744 ± 0.223	0.741 ± 0.223	<b>0.746 ± 0.224</b>
KNN	G-mean	0.827 ± 0.160	0.829 ± 0.158	0.839 ± 0.146	<b>0.840 ± 0.147</b>
LR	Accuracy	0.778 ± 0.157	<b>0.791 ± 0.158</b>	0.764 ± 0.184	0.773 ± 0.185
LR	F-score	0.693 ± 0.243	0.717 ± 0.241	0.718 ± 0.222	<b>0.727 ± 0.226</b>
LR	G-mean	0.796 ± 0.171	0.814 ± 0.165	0.839 ± 0.130	<b>0.842 ± 0.137</b>
RF	Accuracy	0.827 ± 0.145	<b>0.832 ± 0.148</b>	0.827 ± 0.154	0.829 ± 0.153
RF	F-score	0.767 ± 0.215	0.775 ± 0.216	0.781 ± 0.204	<b>0.784 ± 0.204</b>
RF	G-mean	0.844 ± 0.148	0.849 ± 0.149	0.863 ± 0.131	<b>0.865 ± 0.131</b>

TABLE 7

Friedman test results. Statistical significance is tested at a level of  $\alpha = 0.05$ . The null hypothesis is that there is no difference in the classification outcome across oversamplers.

Classifier	Evaluation Metric	p-value	Significance
DT	Accuracy	1.1e-15	True
DT	F-score	2.4e-31	True
DT	G-mean	2.3e-23	True
KNN	Accuracy	5.9e-20	True
KNN	F-score	8.8e-69	True
KNN	G-mean	8.8e-52	True
LR	Accuracy	1.1e-30	True
LR	F-score	4.0e-98	True
LR	G-mean	2.3e-83	True
RF	Accuracy	2.8e-26	True
RF	F-score	1.8e-88	True
RF	G-mean	1.8e-61	True

TABLE 9

Adjusted p-values using the Holm-Bonferroni method. Bold values are statistically significant at a level of  $\alpha = 0.05$ . The null hypothesis is that the Oversampling or Proposed method does not perform better than the control method (Standard AL framework).

Classifier	Evaluation Metric	Oversampling	Proposed
DT	Accuracy	7.7e-01	<b>1.1e-04</b>
DT	F-score	6.3e-02	<b>2.0e-06</b>
DT	G-mean	<b>1.0e-08</b>	<b>2.9e-12</b>
KNN	Accuracy	<b>1.0e-02</b>	8.5e-01
KNN	F-score	<b>7.1e-07</b>	<b>8.3e-13</b>
KNN	G-mean	<b>1.9e-11</b>	<b>1.0e-12</b>
LR	Accuracy	<b>3.2e-02</b>	8.3e-01
LR	F-score	<b>1.5e-09</b>	<b>5.8e-17</b>
LR	G-mean	<b>1.9e-13</b>	<b>5.6e-16</b>
RF	Accuracy	4.3e-01	4.3e-01
RF	F-score	<b>1.4e-11</b>	<b>1.1e-12</b>
RF	G-mean	<b>1.5e-10</b>	<b>1.2e-10</b>

TABLE 8

Adjusted p-values using the Wilcoxon signed-rank method. Bold values are statistically significant at a level of  $\alpha = 0.05$ . The null hypothesis is that the performance of the proposed framework is similar to that of the oversampling or standard framework.

Dataset	Oversampling	Standard
Baseball	5.0e-01	3.4e-01
Gas Drift	<b>3.7e-26</b>	<b>4.6e-57</b>
Gesture Segmentation	<b>1.3e-02</b>	<b>8.7e-04</b>
Image Segmentation	<b>9.6e-18</b>	<b>2.1e-44</b>
Japanese Vowels	<b>2.4e-09</b>	<b>1.6e-32</b>
Mfeat Zernike	<b>1.2e-12</b>	<b>9.5e-40</b>
Mice Protein	<b>6.5e-32</b>	<b>1.5e-61</b>
Pendigits	<b>5.0e-18</b>	<b>2.3e-45</b>
Steel Plates	<b>3.4e-04</b>	<b>1.3e-08</b>
Texture	<b>1.5e-22</b>	<b>6.7e-57</b>
Usps	3.8e-01	<b>2.1e-29</b>
Vehicle	<b>7.4e-11</b>	<b>7.9e-13</b>
Volkert	2.5e-01	<b>1.3e-02</b>
Waveform	<b>8.9e-08</b>	<b>2.6e-02</b>
Wine Quality	<b>3.8e-05</b>	<b>6.1e-03</b>

This study discussed the effect of data augmentation within the AL framework, along with the exploration of optimal augmentation methods within AL iterations. However, the conceptual nature of this study implies some limitations. Specifically, the large number of experiments required to test the method's efficacy, along with the limited computational power available, led to a limited exploration

of the grid search's potential. Future work should focus on understanding how the usage of a more comprehensive parameter tuning approach improves the quality of the AL method. In addition, the proposed method was not able to outperform the standard AL method at 100% of scenarios. The exploration of other, more complex data augmentation techniques might further improve its performance by producing more meaningful training observations. Specifically, in this study, we assume that all datasets used follow a manifold, allowing the usage of G-SMOTE as a data augmentation approach. However, this method cannot be used in more complex, non-euclidean spaces. In this scenario, the usage of G-SMOTE is not valid and might lead to the production of noisy data. Deep Learning-based data augmentation techniques are able to address this limitation and improve the overall quality of the artificial data being generated. We also encountered significant standard errors throughout our experimental results (see Subsection 5.1), consistent with the findings in [20, 46]. This facet suggests that the usage of more robust generators did not decrease the standard error of AL performance. Instead, AL's performance variability is likely dependent on the quality of its initialization.

## 6 CONCLUSION

The ability to train ML classifiers is usually limited to the availability of labeled data. However, manually labeling data is often expensive, which makes the usage of AL particularly appealing for selecting the most informative observations and reducing the amount of required labeled data. On the other hand, the introduction of data variability in the training dataset can also be conducted via data augmentation. However, most, if not all, AL configurations that use some form of data augmentation are domain and/or task-specific. These methods typically apply deep learning approaches to both classification and data augmentation. Consequently, they may not apply to other classification tasks or when the available computational power is insufficient.

In this paper, we proposed a domain-agnostic AL framework that implements Data Augmentation and hyperparameter tuning. We found that a heuristic Data Augmentation algorithm is sufficient to improve the data selection efficiency in AL. Specifically, the data augmentation method used almost always increased AL performance, regardless of the target goal (*i.e.*, optimizing classification or data selection efficiency). The usage of data augmentation reduced the number of iterations required to train a classifier with a performance as good as (or better than) classifiers trained with the entire training dataset (*i.e.*, without using AL). In addition, the proposed method reduced the size of the training dataset, which is expanded with artificial data.

With this revised AL configuration, data selection in AL iterations aims towards observations that optimize the quality of the artificial data produced. The substitution of less informative labeled data with artificial data is especially useful in this context since it reduces some of the user interaction necessary to reach a sufficiently informative dataset. In order to further improve the proposed method, future work should (1) focus on the development of methods with varying data augmentation policies depending on the different input space regions, (2) develop augmentation-sensitive query functions capable of avoiding the unnecessary selection of similar observations from the unlabeled dataset and (3) understand the gap between heuristic/input space data augmentation techniques and neural network/feature space data augmentation techniques in an AL context better.

## ACKNOWLEDGMENTS

This research was supported by three research grants of the Portuguese Foundation for Science and Technology ("Fundação para a Ciência e a Tecnologia"), references SFRH/BD/151473/2021, DSAIPA/DS/0116/2019 and PCIF/SSI/0102/2017.

## REFERENCES

- [1] V. Nath, D. Yang, B. A. Landman, D. Xu, and H. R. Roth, "Diminishing uncertainty within the training pool: Active learning for medical image segmentation," *IEEE Transactions on Medical Imaging*, vol. 40, no. 10, pp. 2534–2547, 2021.
- [2] Y. Sverchkov and M. Craven, "A review of active learning approaches to experimental design for uncovering biological networks," *PLoS Computational Biology*, vol. 13, p. e1005466, 6 2017.
- [3] X. Li, D. Kuang, and C. X. Ling, "Active learning for hierarchical text classification," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 7301 LNAI, pp. 14–25, 2012.
- [4] Y. Li, J. Yin, and L. Chen, "Seal: Semisupervised adversarial active learning on attributed graphs," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 7, pp. 3136–3147, 2021.
- [5] O. Siméoni, M. Budnik, Y. Avrithis, and G. Gravier, "Rethinking deep active learning: Using unlabeled data at model training," *Proceedings - International Conference on Pattern Recognition*, pp. 1220–1227, 2020.
- [6] J. E. Van Engelen and H. H. Hoos, "A survey on semi-supervised learning," *Machine Learning*, vol. 109, no. 2, pp. 373–440, 2020.
- [7] O. Sener and S. Savarese, "Active learning for convolutional neural networks: A core-set approach," in *International Conference on Learning Representations*, 2018.
- [8] Y. Leng, X. Xu, and G. Qi, "Combining active learning and semi-supervised learning to construct svm classifier," *Knowledge-Based Systems*, vol. 44, pp. 121–131, 2013.
- [9] H. Yu, X. Yang, S. Zheng, and C. Sun, "Active learning from imbalanced data: A solution of online weighted extreme learning machine," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, pp. 1088–1103, 4 2019.
- [10] H. Zhang, W. Liu, and Q. Liu, "Reinforcement online active learning ensemble for drifting imbalanced data streams," *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [11] W. Zong, G.-B. Huang, and Y. Chen, "Weighted extreme learning machine for imbalance learning," *Neurocomputing*, vol. 101, pp. 229–242, 2013.
- [12] J. Qin, C. Wang, Q. Zou, Y. Sun, and B. Chen, "Active learning with extreme learning machine for online imbalanced multiclass classification," *Knowledge-Based Systems*, vol. 231, p. 107385, 2021.
- [13] W. Liu, H. Zhang, Z. Ding, Q. Liu, and C. Zhu, "A comprehensive active learning method for multiclass imbalanced data streams with concept drift," *Knowledge-Based Systems*, vol. 215, p. 106778, 2021.
- [14] A. Tharwat and W. Schenck, "Balancing exploration and exploitation: A novel active learner for imbalanced data," *Knowledge-Based Systems*, vol. 210, p. 106500, 2020.
- [15] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Transactions on knowledge and data engineering*, vol. 21, no. 9, pp. 1263–1284, 2009.
- [16] Y.-Y. Kim, K. Song, J. Jang, and I.-c. Moon, "Lada: Look-ahead data acquisition via augmentation for deep active learning," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [17] G. Douzas and F. Bacao, "Geometric SMOTE a geometrically enhanced drop-in replacement for SMOTE," *Information Sciences*, vol. 501, pp. 118–135, Oct. 2019.
- [18] J. Katz-Samuels, J. Zhang, L. Jain, and K. Jamieson, "Improved algorithms for agnostic pool-based active classification," in *International Conference on Machine Learning*. PMLR, 2021, pp. 5334–5344.
- [19] T. Su, S. Zhang, and T. Liu, "Multi-spectral image classification based on an object-based active learning approach," *Remote Sensing*, vol. 12, p. 504, 2 2020.
- [20] J. Fonseca, G. Douzas, and F. Bacao, "Increasing the Effectiveness of Active Learning: Introducing Artificial Data Generation in Active Learning for Land Use/Land Cover Classification," *Remote Sensing 2021, Vol. 13, Page 2619*, vol. 13, no. 13, p. 2619, jul 2021.
- [21] D. Yoo and I. S. Kweon, "Learning loss for active learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 93–102.
- [22] H. H. Aghdam, A. Gonzalez-Garcia, A. Lopez, and 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.
- [23] T. M. Hospedales, S. Gong, and T. Xiang, "Finding rare classes: Active learning with generative and discriminative models," *IEEE transactions on knowledge and data engineering*, vol. 25, no. 2, pp. 374–386, 2011.
- [24] B. Gu, Z. Zhai, C. Deng, and H. Huang, "Efficient active learning by querying discriminative and representative samples and fully exploiting unlabeled data," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, pp. 4111–4122, 9 2021.
- [25] P. Kumar and A. Gupta, "Active learning query strategies for classification, regression, and clustering: A survey," *Journal of Computer Science and Technology* 2020 35:4, vol. 35, pp. 913–945, 7 2020.
- [26] Y. Fu, X. Zhu, and B. Li, "A survey on instance selection for active learning," *Knowledge and information systems*, vol. 35, no. 2, pp. 249–283, 2013.

- [27] S. Behpour, K. M. Kitani, and B. D. Ziebart, "Ada: Adversarial data augmentation for object detection," *Proceedings - 2019 IEEE Winter Conference on Applications of Computer Vision, WACV 2019*, pp. 1243–1252, 3 2019.
- [28] J. Fonseca, G. Douzas, and F. Bacao, "Improving imbalanced land cover classification with k-means smote: Detecting and oversampling distinctive minority spectral signatures," *Information*, vol. 12, no. 7, p. 266, 2021.
- [29] T. DeVries and 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, 2 2017.
- [30] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, no. 1, pp. 1–48, 2019.
- [31] O. Kashefi and R. Hwa, "Quantifying the evaluation of heuristic methods for textual data augmentation," in *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*. Online: Association for Computational Linguistics, Nov. 2020, pp. 200–208.
- [32] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 001–13 008.
- [33] T. Tran, T.-T. Do, I. Reid, and G. Carneiro, "Bayesian generative active deep learning," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6295–6304.
- [34] S. Sinha, S. Ebrahimi, and T. Darrell, "Variational adversarial active learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5972–5981.
- [35] K. Kim, D. Park, K. I. Kim, and S. Y. Chun, "Task-aware variational adversarial active learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 8166–8175.
- [36] Y. Ma, S. Lu, E. Xu, T. Yu, and L. Zhou, "Combining active learning and data augmentation for image classification," in *Proceedings of the 2020 3rd International Conference on Big Data Technologies*, 2020, pp. 58–62.
- [37] H. Quteineh, S. Samothrakis, and R. Sutcliffe, "Textual data augmentation for efficient active learning on tiny datasets," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 7400–7410.
- [38] Q. Li, Z. Huang, Y. Dou, and Z. Zhang, "A framework of data augmentation while active learning for chinese named entity recognition," in *International Conference on Knowledge Science, Engineering and Management*. Springer, 2021, pp. 88–100.
- [39] C. Wu, *The decision tree approach to classification*. Purdue University, 1975.
- [40] T. Cover and P. Hart, "Nearest neighbor pattern classification," *IEEE Transactions on Information Theory*, vol. 13, no. 1, pp. 21–27, 1967.
- [41] T. K. Ho, "Random decision forests," in *Proceedings of the Third International Conference on Document Analysis and Recognition (Volume 1) - Volume 1*, ser. ICDAR '95. USA: IEEE Computer Society, 1995, p. 278.
- [42] J. A. Nelder and R. W. Wedderburn, "Generalized linear models," *Journal of the Royal Statistical Society: Series A (General)*, vol. 135, no. 3, pp. 370–384, 1972.
- [43] L. A. Jeni, J. F. Cohn, and 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.
- [44] M. Fatourehchi, R. K. Ward, S. G. Mason, J. Huggins, A. Schloegl, and G. E. Birch, "Comparison of evaluation metrics in classification applications with imbalanced datasets," in *2008 seventh international conference on machine learning and applications*. IEEE, 2008, pp. 777–782.
- [45] 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.
- [46] D. Kottke, A. Calma, D. Huseljic, G. Kreml, and B. Sick, "Challenges of reliable, realistic and comparable active learning evaluation," in *CEUR Workshop Proceedings*, vol. 1924, sep 2017, pp. 2–14.
- [47] T. Reitmaier and B. Sick, "Let us know your decision: Pool-based active training of a generative classifier with the selection strategy 4ds," *Information Sciences*, vol. 230, pp. 106–131, 5 2013.
- [48] J.-F. Kagy, T. Kayadelen, J. Ma, A. Rostamizadeh, and J. Strnadova, "The practical challenges of active learning: Lessons learned from live experimentation," *arXiv preprint arXiv:1907.00038*, 6 2019.
- [49] 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, and É. Duchesnay, "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [50] G. Lemaître, F. Nogueira, and C. K. Aridas, "Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning," *Journal of Machine Learning Research*, vol. 18, no. 17, pp. 1–5, 2017.
- [51] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," *Journal of Machine Learning Research*, vol. 7, pp. 1–30, 2006.
- [52] M. Friedman, "The use of ranks to avoid the assumption of normality implicit in the analysis of variance," *Journal of the american statistical association*, vol. 32, no. 200, pp. 675–701, 1937.
- [53] F. Wilcoxon, "Individual Comparisons by Ranking Methods," *Biometrics Bulletin*, vol. 1, no. 6, p. 80, dec 1945.
- [54] S. Holm, "A simple sequentially rejective multiple test procedure," *Scandinavian journal of statistics*, pp. 65–70, 1979.



**João Fonseca** is a PhD student at NOVA Information Management School (NOVA IMS) working on data augmentation, active learning and data preprocessing methods. Specifically, I am focused on improving the quality of Land Use/Land Cover classification tasks through the application of these data preprocessing methods. His work is being funded with a MIT Portugal PhD Grant (2020 FCT-MPP2030). In the past, João Fonseca conducted research on Land Use/Land Cover classification methods to automatically update LULC maps of the Portuguese mainland. His research included the development of pipelines to systematize the preprocessing of Sentinel-2 satellite imagery. He also developed and deployed different types of algorithms for various tasks, such as data filtering, dimensionality reduction, feature extraction and classification.



**Fernando Bação** is currently a Full Professor at the NOVA Information Management School (NOVA IMS), Universidade Nova de Lisboa, where he has served as Associate Dean, Director of the MagIC Research Center and Director of the Doctoral Program in Information Management. He holds a PhD in Information Management, from Universidade Nova de Lisboa, and his research interests include data science, machine learning, remote sensing, and information management. His research work has appeared in scientific journals such as *Information Sciences*, *Expert Systems with Applications*, *Information Sciences*, *Computers in Human Behavior*, *The Internet and Higher Education*, *Information & Management*, *International Journal of Geographical Information Science*, among others and received awards for its impact.