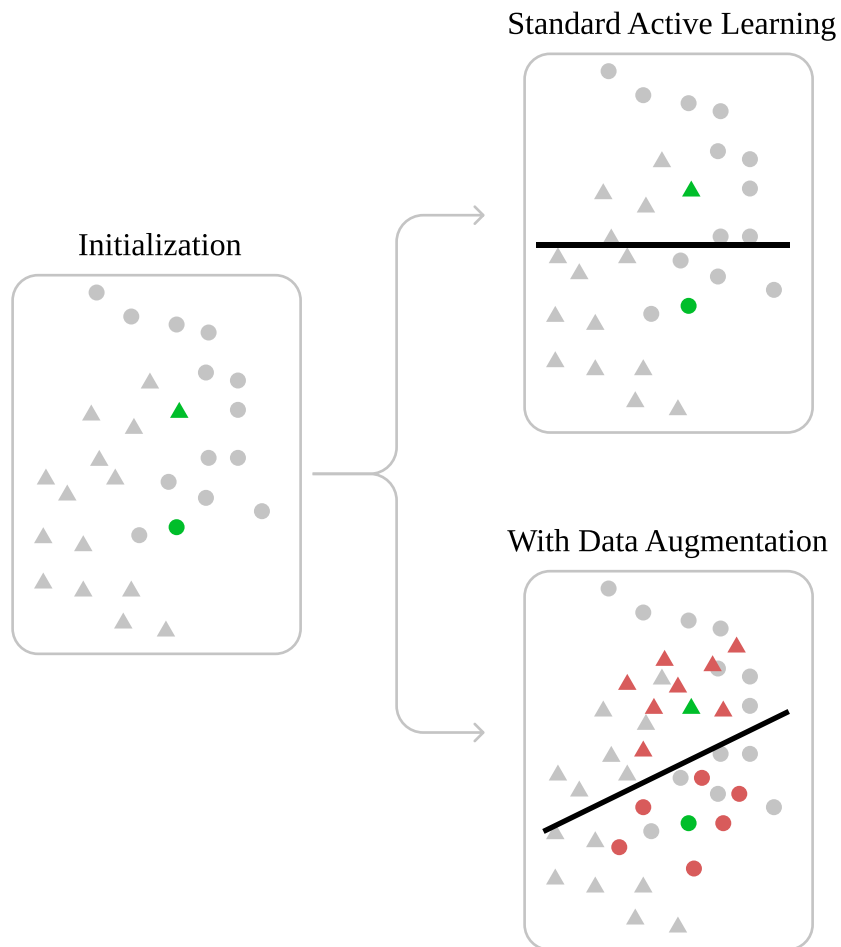


Graphical Abstract

Improving Active Learning Performance Through the Use of Data Augmentation

Joao Fonseca, Fernando Bacao



Highlights

Improving Active Learning Performance Through the Use of Data Augmentation

Joao Fonseca, Fernando Bacao

- We propose a new Active Learning framework that leverages hyperparameter optimization and data augmentation techniques;
- The use of data augmentation in Active Learning is sufficient to substantially improve the performance of an Active Learner, regardless of the choice of dataset/domain, classifier or metric.
- In most scenarios, the proposed method outperformed classifiers trained in fully supervised settings while using less data.

Improving Active Learning Performance Through the Use of Data Augmentation

Joao Fonseca^a, Fernando Bacao^a

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

Abstract

Active Learning (AL) is a technique that is used to iteratively select 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. However, the manual labeling of observations involves human resources with domain expertise, making it an expensive and time-consuming task. The literature describes various methods to improve the effectiveness of this process, but there is little research developed around the usage of artificial data sources in AL. In this paper we propose a new framework for AL, which allows for an effective use of artificial data. Our method implements a data augmentation policy that optimizes the generation of artificial instances to improve the AL process. We compare the proposed method to the standard framework as well as another active learning method that uses data augmentation. The models' performance was tested using 4 different classifiers, 2 AL-specific performance metrics and 3 classification performance metrics over 10 different datasets. We show that the proposed framework, using data augmentation, significantly improves the performance of AL, both in terms of classification performance and data selection efficiency.

Keywords: 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 small amount of the data being produced and stored can be useful for supervised learning tasks. Additionally, it's important 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 special case of ML because it falls between the supervised and unsupervised learning perspectives. AL, instead of optimizing the informativeness of an existing training set, it 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]. 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, 10, 11]. However, classifier or query function-level modifications (See Section 2) have limited applicability since a universally good AL strategy has not been found [7]. Other methods address imbalance learning by weighing the observations as a function of the observation's class imbalance ratio [12]. Alternatively, other methods reduce the imbalanced learning bias by combining Informative and Representative-based query approaches (see Section 2) [13]. Another approach to deal with imbalanced data and data scarcity, in general, is data augmentation. This approach has the advantage of being classifier-agnostic, potentially reduces the imbalanced learning bias, and also works as a regularization method in

data-scarce environments, such as AL implementations [14]. 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 3.1) 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 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 the non-sophisticated user. 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: fixed data augmentation (similar to oversampling strategies, 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 Geometric Synthetic Minority Oversampling Technique (G-SMOTE) [15] as a data augmentation method with an optimized generation policy (explained in Section 3). The hyperparameter optimization module is used to find the best data augmentation policy at each iteration. We test the effectiveness of the proposed method in 10 datasets of different domains. We implement 3 AL frameworks

(standard, constant data augmentation and varying data augmentation) using 4 different classifiers, 3 different performance metrics and calculate 2 AL-specific performance metrics.

The rest of this manuscript is structured as follows: Section 2 describes the state-of-the-art in AL. Section 3 describes the state-of-the-art in Data Augmentation. Section 4 describes the proposed method. Section 5 describes the methodology of the study’s experiment. Section 6 presents the results obtained from the experiment, as well as a discussion of these results. Section 7 presents the conclusions drawn from this study.

2. Active Learning Methods

AL maximizes a classifier’s performance while annotating as least observations as possible. It assumes that observations within the same dataset have a different contribution to the training of ML classifiers [16]. Consequently, the data annotation cost can be reduced with the annotation of the most valuable observations within an unlabeled input space. The goal is to iteratively maximize the classification performance of ML algorithms while reducing the amount of training data required to reach a performance threshold [17]. It allows the implementation of ML classifiers with a good performance and minimal effort when compared to the random selection of data or labeling the entire unlabeled dataset [18]. Therefore, it addresses the labeling problem in scenarios with a limited budget or time.

AL methods are divided in 2 stages, initialization and iteration. Assuming the AL task is initialized without any previously labeled data, it is typically collected randomly across an unlabeled pool of data, which is then labeled by a supervisor [19, 20, 21].

Once an initial training dataset is set up, the iterative process of AL takes place. An AL iteration is completed once a new batch of labeled data is added to the training dataset. A standard AL process is shown in Figure 1 and is composed of the following steps [22, 2]:

1. Setting up a classification algorithm and uncertainty criterion. The classifier is trained using the labeled dataset (*i.e.*, the Current Training Dataset), and is used to predict the class membership probabilities of the observations found in the unlabeled dataset. The class probabilities are passed into an Uncertainty Criterion, which will return the

116 classification uncertainty of the classification algorithm for each un-
 117 labeled observation. The combination of the classifier, along with the
 118 uncertainty criterion is sometimes referred to as the Query/Acquisition
 119 function [23].

120 2. Selecting the top N observations. Since it is not possible to determine a
 121 priori whether the classifier’s prediction is correct or not, the N obser-
 122 vations with highest uncertainty may have been unknowingly correctly
 123 classified. However, regardless of the classification quality, these obser-
 124 vations are expected to provide the most meaningful information to
 125 train the classifier in the next iteration.

126 3. Labeling the selected N observations and updating the current training
 127 dataset with the new training observations. The selected observations
 128 from the unlabeled dataset are presented to the supervisor, which is
 129 responsible for manually labeling the observations. The new (labeled)
 130 training observations are added to the training dataset and the iteration
 131 is completed.

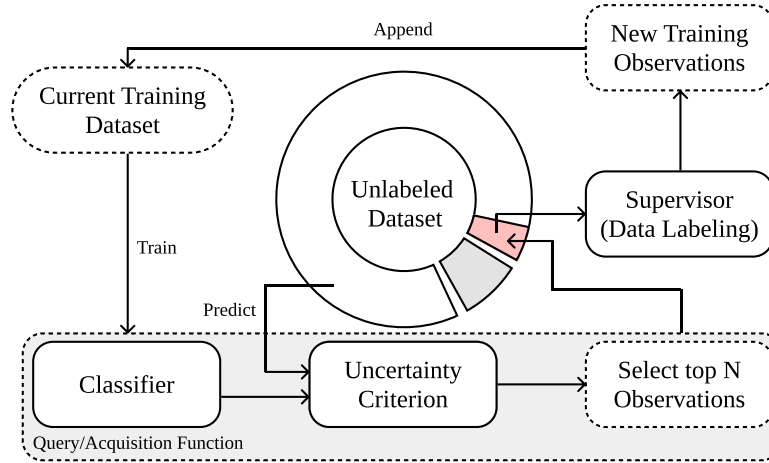


Figure 1: Diagram depicting a typical AL iteration. In the first iteration, the training set collected during the initialization process becomes the “Current Training Dataset”.

132 Two common challenges found in AL implementations is the consistency
 133 and efficiency of AL in practical scenarios [24]. On the one hand, the consi-

134 tency problem refers to the high variance in performance (regarding classifi-
 135 cation and data selection) over different initializations (*i.e.*, different initial
 136 training datasets) of active learners. On the other hand, the efficiency prob-
 137 lem refers to the maximization of the quality of the collected data over a run.
 138 Therefore, a good active learner is capable of having a consistent performance
 139 over different initializations while ensuring the production of high-performing
 140 classifiers with the least possible amount of data. There are various factors
 141 that may affect the consistency and efficiency of the AL framework: (1)
 142 Human error during data labeling [25], (2) Non-informative initial training
 143 dataset [26] and (3) Lack of an appropriate uncertainty criterion [23]. AL re-
 144 search has typically been focused on the specification of uncertainty criteria,
 145 as well as domain-specific applications. Query functions can be divided into
 146 two different categories [27, 28]:

- 147 1. Informative-based query strategies. These strategies use the classifier’s
 148 output to assess the importance of each observation towards the perfor-
 149 mance of the classifier. These strategies focus on quantifying the class
 150 uncertainty of the unlabeled observations. Since these techniques do
 151 not account for the relationships between the unlabeled observations
 152 and treats each observation independently [29].
- 153 2. Representative-based query strategies. These strategies estimate the
 154 optimal set of observations that will optimize the classifier’s perfor-
 155 mance. This strategy contains 3 main approaches: Density-based,
 156 Diversity-based and Exploration of graph structures. Although this
 157 method addresses the problem of sampling bias and redundant instance
 158 selection, these strategies typically require more observations in order
 159 to reach the desired classification performance [28].

160 Although there are significant contributions towards the development of
 161 more robust query functions and classifiers in AL, modifications to AL’s
 162 basic structure is rarely explored. In [20] the authors introduce a loss predic-
 163 tion module in the AL framework to replace the uncertainty criterion. This
 164 model implements a second classifier to predict the expected loss of the un-
 165 labeled observations (using the actual losses collected during the training of
 166 the original classifier) and return the unlabeled observations with the highest
 167 expected loss. Although this contribution is specific to neural networks (and
 168 more specifically, to deep neural networks), they were able to significantly
 169 improve the efficiency of data selection in AL. In [5] the authors propose the

usage of semi-supervised learning during both the initialization of the AL and the iterative process as well. However, this method was proposed specifically for deep learning applications. In [19], the authors introduce the generator element in the AL framework (discussed in Section 4) using an oversampling method, showing that this method effectively addresses the limitations of imbalanced learning. However, this method was implemented specifically in the Remote Sensing domain and used an oversampling strategy without consideration for the actual amount of artificial data generated, which may limit its performance.

2.1. Query Strategies

A query strategy/function encompasses all the steps prior to the data labeling within an AL iteration. They focus on finding the observations' informativeness, representativeness or both [27, 28]. Representative query strategies are generally less efficient in data selection than Informative query strategies [28]. However, recent research often use representative approaches alongside informative approaches [27, 30]. Representative query strategies are explored via 3 main approaches [28]:

1. Density-based, which select representative observations from high density regions. [31, 3, 32] used a density-based approach using clustering algorithms to select the observations closest to the centroid of each cluster.
2. Diversity-based, which select the N observations at each iteration that maximize the diversity in the training data. The diversity-based approach was developed to avoid the selection of redundant observations in batch-mode learning [33].
3. Graph-based, which find the most representative nodes and edges of a graph network [34]. Since these methods are specific to graph network data, they have a more limited applicability.

Informative query strategies, unlike representative query strategies, do not account for the structure of the unlabeled dataset. As a result, this type of strategy may lead to the inefficient selection of observations (*i.e.*, redundant observations with similar profiles) [28]. Research on more robust selection criteria attempts to address the efficiency problem. This is motivated by

the importance of the selection criteria in AL’s iterative process [23]. Specifically, Settles [35] observed that in some datasets informative query strategies fail to outperform the random selection of observations. Generally, the Random Selection query method is used as a baseline. This method disregards the class membership probabilities produced by the classifier and returns N random points from the dataset without following any specific criteria.

A frequently used query strategy is Uncertainty Sampling, originally proposed in [36]. Using this method, the estimation of an observation’s uncertainty is based on the target class with the highest probability (p_a , according to the classifier) and the uncertainty is calculated as $1 - p_a$. However, since this method dismissed the classifier’s predictions on the remaining labels, the Breaking Ties criterion was proposed to address this limitation for multiclass problems [37]. This method uses the two target classes with highest probability (p_a and p_b , according to the classifier) and the uncertainty is calculated as $p_a - p_b$ (in this case, the lower the output value, the higher the uncertainty). Recent variants of the Breaking Ties criterion, such as the Modified Breaking Ties, attempted to fix some limitations of the original method [38, 39].

Another common informative query strategy is the calculation of Shannon’s Entropy. This metric measures the level of uncertainty based on the probabilities of a set of possible events. Its formula is given by $H(p) = -\sum_{i=0}^n p_i \log_2 p_i$, having p as the set of probabilities of all target classes. The application of the Entropy uncertainty criterion is also frequently applied in Deep Active Learning [21]. Other Entropy-based methods were also developed for more specific applications. For example, an ensemble querying approach known as Entropy Querying-by-Bagging uses the predictions of all estimators to find the maximum entropy of each observation [40].

The Query by Committee (QBC) strategy was developed to address ensemble classifiers. It is a disagreement based strategy attempts to maximize the information gain at each iteration by computing the disagreement of the predictions over the estimators that form the ensemble. The Entropy Querying-by-Bagging and Query-by-Boosting methods are also ensemble strategies. Query by boosting and bagging methods were found to achieve a good performance over various datasets [41], while the performance between the two strategies appears to differ significantly across various scenarios [42].

Other classifier-specific query strategies were also developed for different applications. However, these methods have the disadvantage of depending on the classifier being used. For example, Margin Sampling is a well studied strategy that uses a Support Vector Machine as its classifier in order to select

241 the unlabeled observations closest to its decision boundaries [28]. Although,
242 since this method is known to lead to the excessive selection of observations
243 in dense regions [43], it was improved in various ways. In [43] the authors
244 extend this strategy by applying the manifold-preserving graph reduction
245 algorithm beyond the normal Margin Sampling method.

246 **3. Data Augmentation Methods**

247
248 Data Augmentation methods expand the training dataset by introduc-
249 ing new and informative observations [44]. The production of artificial data
250 may be done via the introduction of perturbations on the input [45], fea-
251 ture [46] or output space [44]. Data Augmentation methods may be divided
252 into Heuristic and Neural Network-based approaches [47]. In addition, they
253 may also be distinguished based on its data generation policy, whether lo-
254 cal (considers a local/specific subset of the dataset) or global (considers the
255 overall distribution of the training dataset). Figure 2 shows the general tax-
256 onomy of Heuristic Data Augmentation methods. Finding the appropriate
257 Data Augmentation method generally depends on the domain [46], although
258 some studies discuss which methods are more appropriate according to the
259 domain [47, 48, 49].

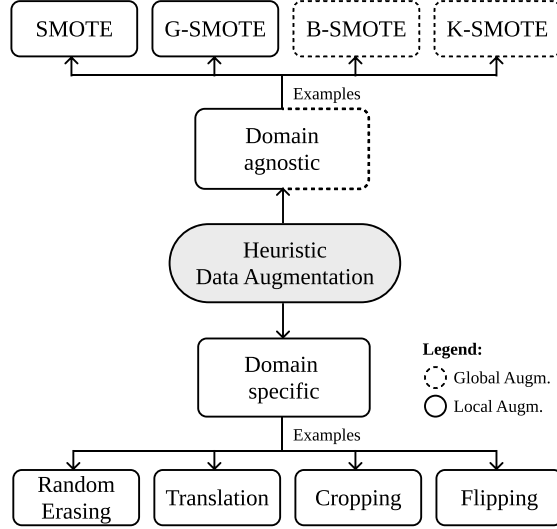


Figure 2: Schema containing a general Heuristic Data Augmentation taxonomy.

Heuristic approaches attempt to generate new and relevant observations through the application a predefined procedure, usually incorporating some degree of randomness [50]. Since these methods typically occur in the input space, they require less data and computational power when compared to Neural Network methods. Neural Network approaches, on the other hand, map the original input space into a lower-dimensional representation, known as the feature space [46]. 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. Considering the scope of this paper (the paper’s contribution is described in Sections 1 and 4), the computational power available for this experiment and the breadth of datasets used in our experimental procedure, we will focus on domain-agnostic heuristic data augmentation methods.

While some techniques may depend on the domain, others are domain-agnostic. For example, Random Erasing [45], Translation, Cropping and Flipping are image data-specific augmentation methods. Other methods, such as most of the variants of the Synthetic Minority Oversampling TEchnique (SMOTE) [51], may be considered domain agnostic. However, SMOTE methods were originally developed as oversamplers, whose goal is to balance

280 the class frequencies of the target variable in the training dataset and ad-
 281 dress the class imbalance bias [52]. Therefore, oversampling methods may be
 282 considered a subset of Data Augmentation. Data Augmentation strategies
 283 may follow varying augmentation strategies, which does not necessarily de-
 284 pend on the target class distribution. An example of the differences among
 285 general data augmentation and oversampling generation strategies is shown
 286 in Figure 3.

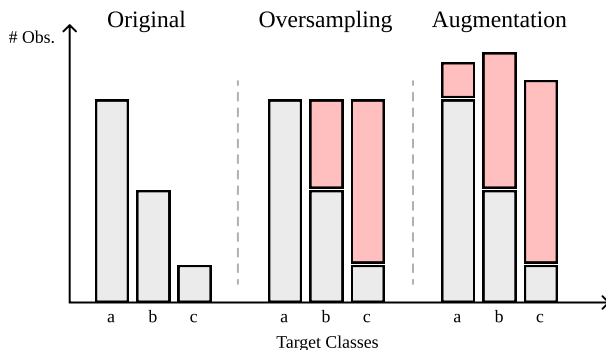


Figure 3: Examples of data augmentation Strategies. The salmon-colored bars represent artificial data using the normal oversampling (center group) and an example of augmentation (right group) strategies.

287 The simplest approach found in the literature is randomly duplicating
 288 existing training observations. As a non-informed data generation method,
 289 although simple to implement, it increases the risk of overfitting and generally
 290 performs worse than other informed heuristic methods [53].

291 The SMOTE method generates artificial data via the linear interpolation
 292 between a random observation and one of its k -nearest neighbors (also ran-
 293 domly selected) [51]. Although simple and effective, it also contains several
 294 limitations which motivated the development other variants, discussed below.
 295 Specifically, its selection mechanism does not consider the global structure of
 296 the dataset while its generation mechanism introduces little variability into
 297 the training dataset [15]. Borderline-SMOTE (B-SMOTE) [54] improves the
 298 selection mechanism by attributing a larger importance to the observations
 299 closer to the decision boundaries. The selected observations are used to run
 300 the SMOTE method in order to produce better defined decision boundaries.
 301 A more recent improvement of the selection mechanism is K-means SMOTE

302 (K-SMOTE) [55]. This method uses a clustering-based approach to over-
 303 come imbalances between and within classes, while considering the densities
 304 of each region of the input space.

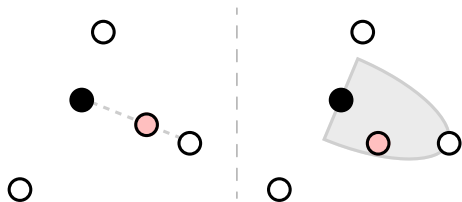


Figure 4: Examples of data generation using SMOTE and G-SMOTE. In this example, both G-SMOTE’s deformation and truncation parameters assume values around 0.5.

305 G-SMOTE [15] modifies SMOTE’s generation mechanism. Instead of
 306 generating an observation as a linear combination between 2 others, it gen-
 307 erates observations within an hypersphere defined using the selected obser-
 308 vation as its center and one of its nearest neighbors as its boundary. The
 309 hypersphere contains two hyperparameters, the truncation and deformation
 310 factors, which limit the area of the hypersphere. The difference between
 311 SMOTE and G-SMOTE is shown in Figure 4. Reference [53] found that
 312 G-SMOTE outperforms various state-of-the-art oversamplers.

313 3.1. Data Augmentation in Active Learning

314
 315 As found in Section 2, the improvement of the AL framework found in the
 316 literature are mostly focused on modifications of the classifier or query strat-
 317 egy. However, a few recent AL applications implementing data augmentation
 318 were found. The method proposed in [14], Look-Ahead Data Acquisition for
 319 Deep Active Learning, implement image data specific data augmentation
 320 to train a deep learning classifier. However, in this study, data augmenta-
 321 tion is based on the unlabeled observations and occurs before the unlabeled
 322 data selection. In [56] the proposed AL method was designed specifically
 323 for image data classification, where a deep learning model was implemented
 324 as a classifier, but its architecture is not described. Other AL frameworks
 325 implementing data augmentation may also be found for Natural Language
 326 Processing applications [57, 58]. However, these methods were designed for

327 specific within that domain and are not necessarily transferrable to other
328 domains or tasks.

329 4. Proposed Method

330
331 Based on the literature found on AL, most of the contributions and
332 novel implementations of AL algorithms focused on the improvement of the
333 choice/architecture of the classifier or the improvement of the uncertainty
334 criterion. In addition, the resulting classification performance of AL-trained
335 classifiers is frequently inconsistent and marginally improve the classifica-
336 tion performance when compared to classifiers trained over the full training
337 set. Finally, in [19] the authors also found a significant variability of the
338 data selection efficiency during different runs of the AL iterative process. In
339 that study the authors proposed a new element within the AL framework,
340 the generator, which was able to marginally reduce the variability previously
341 identified. However, this modification was applied in a Land Use/Land Cover
342 context which contains specific characteristics that are not necessarily found
343 in other supervised learning problems. Specifically, these types of datasets
344 are high dimensional and have limited data variability within each class (*i.e.*,
345 cohesive spectral signatures within classes) due to their geographical prox-
346 imity. Furthermore, the implementation of the generator was done using a
347 simple oversampling augmentation policy, which limits the possibility of em-
348 ploying other techniques with an undefined target amount of data generated
349 at each iteration.

350 This paper provides a context-agnostic AL framework towards the inte-
351 gration of Data Augmentation within AL, with the following contributions:

- 352 1. Improvement of the AL framework by introducing a parameter tuning
353 stage using only the labeled dataset available at the current iteration
354 (*i.e.*, no labeled hold-out set is needed).
- 355 2. Generalization of the generator module proposed in [19] from oversam-
356 pling techniques to any other data augmentation mechanism and/or
357 policy.
- 358 3. Implementation of data augmentation outside of the Deep AL realm,
359 which was not previously found in the literature.

360 4. Analysis of the impact of Data Augmentation and Oversampling in AL
 361 over 10 different datasets of different domains, while comparing them
 362 with the standard AL framework.

363 The proposed iterative process of the AL framework is depicted in Fig-
 364 ure 5. The generator element becomes an additional source of data and is
 365 expected to introduce additional data variability into the training dataset.
 366 This should allow the classifier to generalize better and perform more con-
 367 sistently over unseen observations. However, in this scenario, the amount
 368 of data to generate per class at each iteration is unknown. Consequently,
 369 the hyperparameter tuning step was introduced to estimate the optimal data
 370 augmentation policy at each iteration. In our implementation, this step uses
 371 the current training dataset to perform an exhaustive search over specified
 372 parameters of the generator, tested over a 5-fold cross validation method.
 373 The best augmentation policy found is used to train the iteration’s classifier
 374 in the following step.

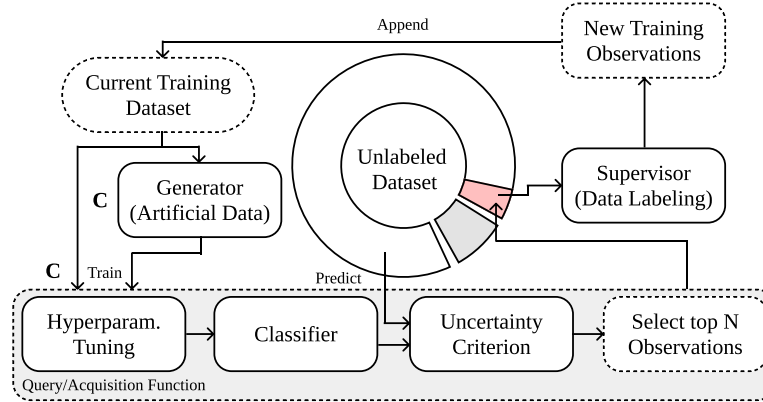


Figure 5: Diagram depicting the proposed AL iteration. The proposed modifications are marked with a boldface “C”.

375 To show the effectiveness of data augmentation in an AL implementation,
 376 we implemented a simple modification of the G-SMOTE algorithm. This
 377 modification facilitates the usage of G-SMOTE beyond its original oversam-
 378 pling purposes. In this paper, the data augmentation strategies used ensure
 379 the all the class frequencies are balanced. Furthermore, the amount of ar-
 380 tificial data produced for each class is defined by the *augmentation factor*,

381 which represents a percentage of the majority class C_{maj} (*e.g.*, an augmenta-
 382 tion factor of 1.2 will ensure there are $\text{count}(C_{maj}) \times 1.2$ observations in every
 383 class). In this paper’s experiment, the data generation mechanism is similar
 384 to the one in [19]. This allows the direct comparison of the two frameworks
 385 and establish a causality of the performance variations to the data generation
 386 mechanism (*i.e.*, augmentation vs normal oversampling) and hyperparameter
 387 tuning steps.

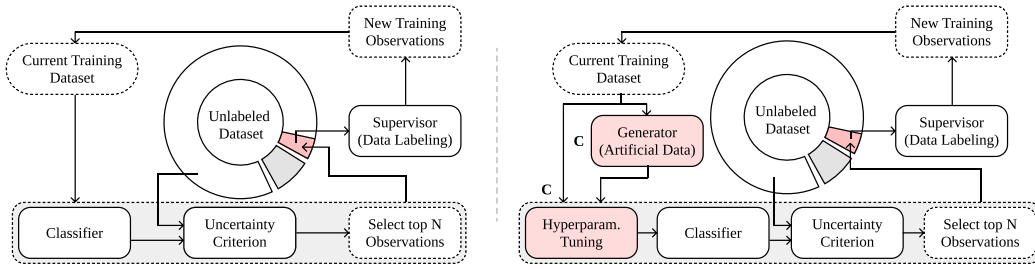


Figure 6: Simplified diagrams highlighting the differences between the proposed and standard AL iterations. The proposed modifications are highlight in red and marked with a boldface “C”.

388 The comparison of diagrams between the proposed and standard AL
 389 frameworks is shown in Figure 6. In the proposed framework, we (1) gen-
 390 eralize the generator module to accept any data augmentation method or
 391 policy and (2) a hyperparameter tuning module to estimate the optimal
 392 data augmentation policy. This framework was designed to be task-agnostic.
 393 Specifically, any data augmentation method (domain specific or not) may be
 394 used, as well as any other parameter search method. It is also expected to
 395 be compatible with other AL modifications, including the ones that do not
 396 affect solely the classifier or uncertainty criterion, such as the one proposed
 397 in [20].

398 5. Methodology

399
 400 This section describes the different elements included in the experimental
 401 procedure. The datasets used were acquired in open data repositories and its
 402 sources and preprocessing steps are defined in Subsection 5.1. The choice of
 403 classifiers used in the experiment are defined in Subsection 5.2. The metrics

404 chosen to measure AL performance and overall classification performance
 405 are defined in Subsection 5.3. The experimental procedure is described in
 406 Subsection 5.4. The implementation of the experiment and resources used
 407 to do so are described in Subsection 5.5.

408 The methodology developed serves 2 purposes: (1) Compare classification
 409 performance once all the AL procedures are completed (*i.e.*, optimal perfor-
 410 mance of a classifier trained via iterative data selection) and (2) Compare the
 411 amount of data required to reach specific performance thresholds (*i.e.*, num-
 412 ber of AL iterations required to reach similar classification performances).

413 5.1. Datasets

414
 415 The datasets used to test the proposed method are publicly available in
 416 open data repositories. Specifically, they were retrieved from OpenML and
 417 the UCI Machine Learning Repository. They were chosen considering dif-
 418 ferent domains of application, imbalance ratios, dimensionality and number
 419 of target classes, all of them focused on classification tasks. The goal is to
 420 demonstrate the performance of the different AL frameworks in various sce-
 421 narios and domains. The data preprocessing approach was similar across all
 422 datasets. Table 1 describes the key properties of the 10 preprocessed datasets
 423 where the experimental procedure was applied.

Dataset	Features	Instances	Minority instances	Majority instances	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
Baseball	15	1320	57	1196	20.98	3

Table 1: Description of the datasets collected after data preprocessing. The sampling strategy is similar across datasets. Legend: (IR) Imbalance Ratio

424 The data preprocessing pipeline is depicted as a flowchart in Figure 7.

425 The missing values are removed from each dataset by removing the corre-
 426 sponding observations. This ensures that the input data in the experiment
 427 is kept as close to its original form as possible. The non-metric features (*i.e.*,
 428 binary, categorical and ordinal variables) were removed since the application
 429 of G-SMOTE is limited to continuous and discrete features. The datasets
 430 containing over 2000 observations were downsampled in order to maintain
 431 the datasets to a manageable size. The data sampling procedure preserves
 432 the relative class frequency of the dataset, in order to maintain the Imbal-
 433 ance Ratio (IR) originally found in each dataset (where $IR = \frac{count(C_{maj})}{count(C_{min})}$).
 434 The remaining features of each dataset are scaled to the range of $[-1, 1]$ to
 435 ensure a common range across features.

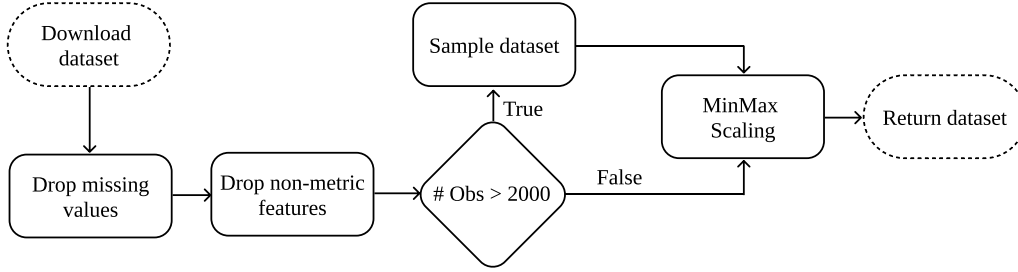


Figure 7: Data preprocessing pipeline.

436 The preprocessed datasets were stored into a SQLite database file and is
 437 available along with the experiment’s source code in the GitHub repository
 438 of the project (see Subsection 5.5).

439 5.2. Machine Learning Algorithms

440
 441 We used a total of 4 classification algorithms and a heuristic data aug-
 442 mentation mechanism. The choice of classifiers was based on the popularity
 443 and family of the classifiers (tree-based, nearest neighbors-based, ensemble-
 444 based and linear models). Our proposed method was tested using a Decision
 445 Tree (DT) [59], a K-nearest neighbors classifier (KNN) [60], a Random For-
 446 est Classifier (RF) [61] and a Logistic Regression (LR) [62]. Since the target
 447 variables are multi-class, the LR classifier was implemented using the one-
 448 versus-all approach. The predicted class is assigned to the label with the
 449 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 4.

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 normal oversampling method, as proposed in [19]. 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 4.

5.3. 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 quantify a model’s classification performance [63]. 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 [64]. However, these metrics remain popular choices for the evaluation of classification performance. Other performance metrics like $Precision = \frac{TP}{TP+FN}$, $Recall = \frac{TP}{TP+FN}$ 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 at a per-class basis instead. In a multiple dataset with varying amount of target classes and meanings, comparing the performance of different models using these metrics becomes impractical.

Based on the recommendations found in [63, 65], we used 2 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 [65]. Both metrics are calculated in a multiclass

486 context considering a one-versus-all approach. For multiclass problems,
 487 the G-mean scorer is calculated as its average per class values:

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

- 488 • The F-score metric consists of the harmonic mean of Precision and
 489 Recall. The two metrics are also calculated considering a one-versus-
 490 all approach. The F-score for the multi-class case can be calculated
 491 using its average per class values [63]:

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

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

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

495 The comparison of the performance of AL frameworks is based on its
 496 data selection and augmentation efficacy. Specifically, an efficient data se-
 497 lection/generation policy allows the production of classifiers with high per-
 498 formance on unseen data while using as least non-artificial training data as
 499 possible. To measure the performance of the different AL setups, we follow
 500 the recommendations found in [24]. The performance of an AL setup will be
 501 compared using two AL-specific performance metrics:

- 502 • Area Under the Learning Curve (AULC). It is the sum of the classifi-
 503 cation performance over a validation/test set of the classifiers trained
 504 of all AL iterations. To facilitate the interpretability of this metric,
 505 the resulting AULC scores are fixed within the range $[0, 1]$ by dividing
 506 the AULC scores by the total amount of iterations (*i.e.*, the maximum
 507 performance area).
- 508 • Data Utilization Rate (DUR) [66]. Measures the percentage of training
 509 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.

5.4. 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 [67]. In this scenario, since the dataset is already labeled, the annotation process is done at zero cost. Figure 8 depicts the experiment designed for one dataset over a single run.

A single run starts with the splitting of a preprocessed dataset in 5 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 5 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 Train set. The AL simulations and the classifiers’ training occur within the Train set. However, the classifiers used to find the maximum performance classification scores are trained over the full Train 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 Train 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 8 is repeated over 3 runs using different random seeds over the 10 different datasets collected. The final scores of each AL configuration and classifier correspond to the average of the 3 runs and 5-fold Cross Validation estimations (*i.e.*, the mean score of 15 fits, across 10 datasets).

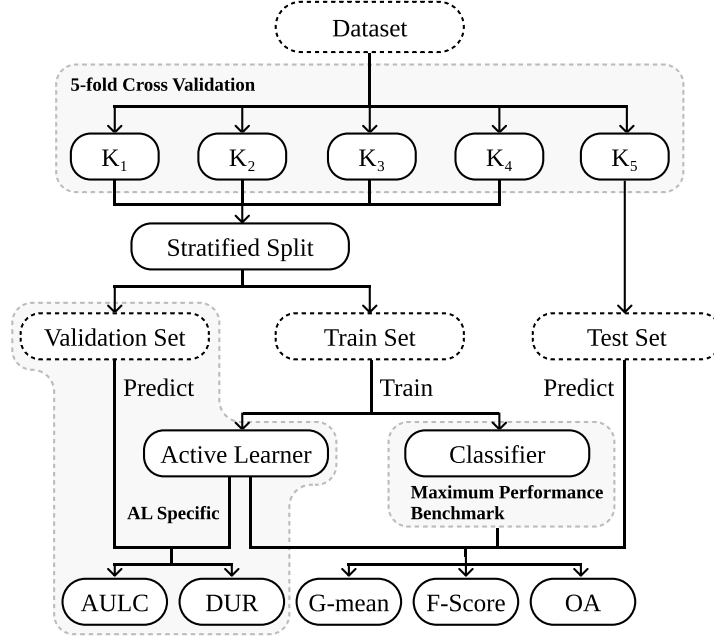


Figure 8: Experimental procedure flowchart. The preprocessed datasets are split into five folds. One of the folds is used to test the best found classifiers using AL and the classifiers trained using the entire training dataset (containing the remaining folds). The training set is used to run both the 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 subsets for overall classification performance and AL-specific performance to avoid data leakage.

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

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

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

554 5.5. *Software Implementation*

555

556 The experiment was implemented using the Python programming lan-
557 guage, along with the Python libraries Scikit-Learn [68], Imbalanced-Learn [69],
558 Geometric-SMOTE [15], Research-Learn and ML-Research libraries. All
559 functions, algorithms, experiments and results are provided in the GitHub
560 repository of the project.

561 6. **Results & Discussion**

562

563 In a multiple dataset experiment, the analysis of results should not rely
564 uniquely on the average performance scores across datasets. The domain of
565 application and fluctuations of performance scores between datasets make
566 the analysis of these averaged results less accurate. Instead, it is generally
567 recommended the use of the mean ranking scores to extend the analysis [70].
568 Since mean performance scores are still intuitive to interpret, we will present
569 and discuss both results. The rank values are assigned based on the mean
570 scores of 3 different runs of 5-fold Cross Validation (15 performance estima-
571 tions per dataset) for each combination of dataset, AL configuration, classifier
572 and performance metric.

573 6.1. *Results*

574

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

Classifier	Evaluation Metric	Standard	Oversampling	Proposed
DT	Accuracy	2.50 ± 0.81	2.20 ± 0.40	1.30 ± 0.64
DT	F-score	2.50 ± 0.81	2.10 ± 0.30	1.40 ± 0.80
DT	G-mean	2.70 ± 0.64	2.00 ± 0.45	1.30 ± 0.64
KNN	Accuracy	2.40 ± 0.80	1.90 ± 0.54	1.70 ± 0.90
KNN	F-score	2.60 ± 0.66	1.80 ± 0.40	1.60 ± 0.92
KNN	G-mean	2.80 ± 0.40	1.70 ± 0.46	1.50 ± 0.81
LR	Accuracy	2.60 ± 0.66	2.10 ± 0.54	1.30 ± 0.64
LR	F-score	2.80 ± 0.40	2.00 ± 0.45	1.20 ± 0.60
LR	G-mean	2.80 ± 0.40	2.00 ± 0.45	1.20 ± 0.60
RF	Accuracy	2.60 ± 0.66	1.90 ± 0.54	1.50 ± 0.81
RF	F-score	2.60 ± 0.66	2.00 ± 0.45	1.40 ± 0.80
RF	G-mean	2.80 ± 0.40	1.60 ± 0.49	1.60 ± 0.80

Table 3: Mean rankings of the AULC metric over the different datasets (10), folds (5) and runs (3) used in the experiment. The proposed method always improves the results of the original framework and on average almost always improves the results of the oversampling framework.

Table 4 shows the average AULC scores, grouped by classifier, Evaluation Metric and AL framework. The variation in performance across active learners is consistent with the mean rankings found in Table 3, while showing significant AULC score differences between the proposed AL method and 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 only the thresholds ending with 0 or 6. 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 9. 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 69.6% 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 30.4% less data).

The mean optimal classification scores of AL methods and Classifiers

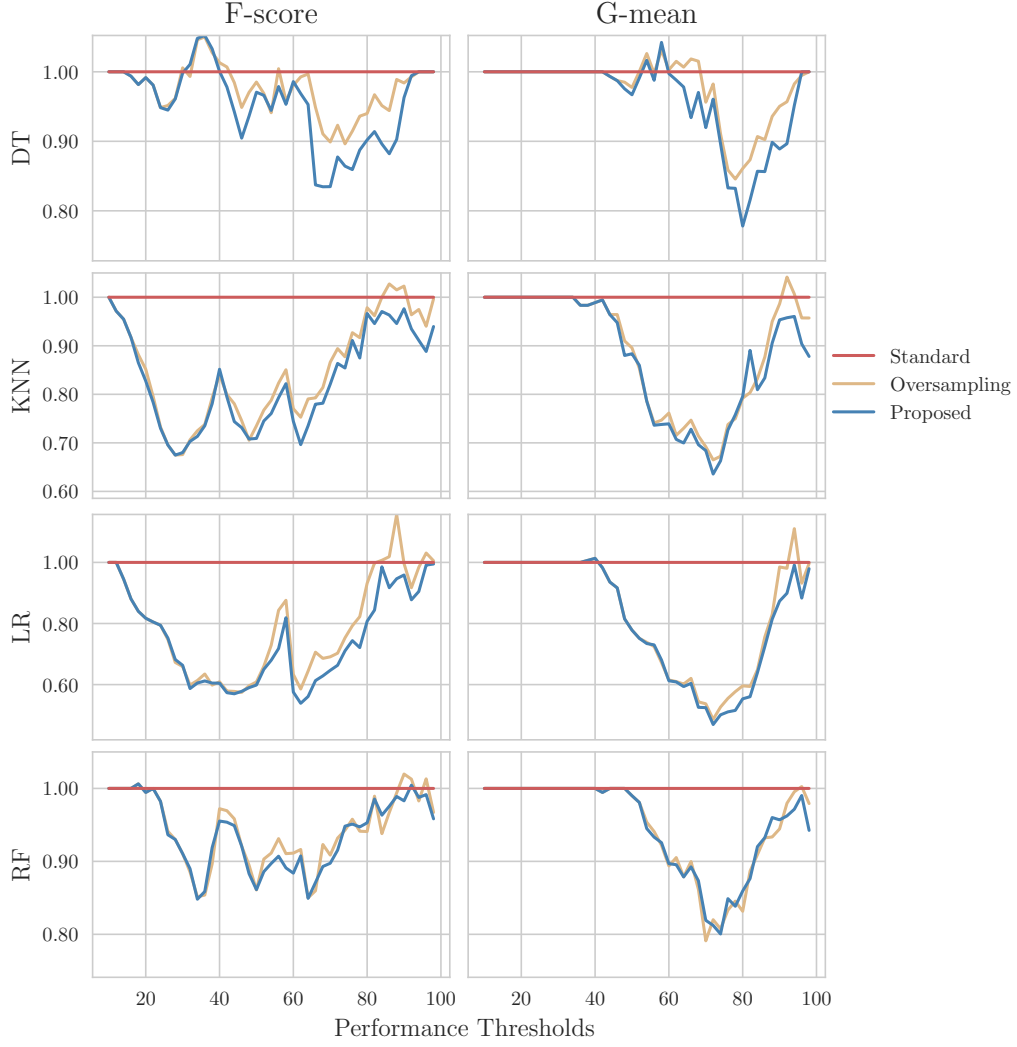


Figure 9: Mean data utilization rates. The y-axis shows the percentage of data (relative to the baseline AL framework) required to reach the different performance thresholds.

Classifier	Evaluation Metric	Standard	Oversampling	Proposed
DT	Accuracy	0.733 ± 0.092	0.732 ± 0.087	0.740 ± 0.087
DT	F-score	0.695 ± 0.088	0.698 ± 0.090	0.705 ± 0.092
DT	G-mean	0.804 ± 0.065	0.811 ± 0.060	0.816 ± 0.062
KNN	Accuracy	0.816 ± 0.091	0.818 ± 0.088	0.822 ± 0.091
KNN	F-score	0.775 ± 0.102	0.784 ± 0.108	0.788 ± 0.111
KNN	G-mean	0.852 ± 0.084	0.866 ± 0.072	0.869 ± 0.074
LR	Accuracy	0.802 ± 0.091	0.812 ± 0.088	0.821 ± 0.086
LR	F-score	0.749 ± 0.112	0.773 ± 0.116	0.784 ± 0.115
LR	G-mean	0.839 ± 0.093	0.870 ± 0.065	0.875 ± 0.064
RF	Accuracy	0.861 ± 0.076	0.861 ± 0.075	0.862 ± 0.077
RF	F-score	0.823 ± 0.105	0.827 ± 0.105	0.829 ± 0.105
RF	G-mean	0.886 ± 0.077	0.895 ± 0.063	0.895 ± 0.065

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).

(fully labeled training set, without AL) is shown in Table 6. The proposed AL method produces classifiers that are almost always able to outperform classifiers using the full training set (*i.e.*, the ones labeled as MP).

6.2. Statistical Analysis

When checking for statistical significance in a multiple dataset context it is important to account for the multiple comparison problem. Consequently, our statistical analysis focuses on the recommendations found in [70]. Overall, we perform 3 statistical tests. The Friedman test [71] is used to understand whether there is a statistically significant difference in performance between the 3 AL frameworks. As post hoc analysis, the Wilcoxon signed-rank test [72] was used 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 analysis, the Holm-Bonferroni [73] method was used to check for statistical significance between the methods using data generators and the Standard AL framework across classifiers and evaluation metrics.

G-mean Score	Classifier	Standard	Oversampling	Proposed
0.60	DT	3.2%	3.1%	3.2%
0.60	KNN	3.6%	2.6%	2.5%
0.60	LR	3.9%	2.2%	2.2%
0.60	RF	2.4%	2.1%	2.1%
0.66	DT	4.6%	4.6%	4.2%
0.66	KNN	4.9%	3.7%	3.5%
0.66	LR	5.7%	3.2%	3.1%
0.66	RF	3.0%	2.8%	2.7%
0.70	DT	6.6%	6.1%	5.8%
0.70	KNN	8.5%	5.0%	4.7%
0.70	LR	9.5%	4.6%	4.3%
0.70	RF	4.5%	3.2%	3.3%
0.76	DT	16.5%	13.0%	12.7%
0.76	KNN	17.8%	9.7%	9.0%
0.76	LR	16.6%	10.0%	7.8%
0.76	RF	10.1%	5.5%	5.5%
0.80	DT	36.1%	30.4%	27.1%
0.80	KNN	22.7%	18.0%	17.8%
0.80	LR	25.2%	16.0%	14.2%
0.80	RF	15.5%	9.0%	9.5%
0.86	DT	60.5%	56.7%	54.5%
0.86	KNN	39.9%	37.0%	37.8%
0.86	LR	32.6%	27.5%	27.0%
0.86	RF	28.0%	25.7%	25.7%
0.90	DT	72.5%	70.7%	67.8%
0.90	KNN	49.9%	50.3%	49.3%
0.90	LR	52.5%	53.8%	49.3%
0.90	RF	44.6%	42.6%	43.5%
0.96	DT	100.0%	99.5%	100.0%
0.96	KNN	79.4%	75.6%	71.6%
0.96	LR	87.5%	83.1%	79.8%
0.96	RF	63.6%	64.2%	63.1%

Table 5: Mean data utilization of AL algorithms, as a percentage of the training set.

Classifier	Evaluation Metric	MP	Standard	Oversampling	Proposed
DT	Accuracy	0.809 ± 0.086	0.802 ± 0.089	0.806 ± 0.089	0.812 ± 0.087
DT	F-score	0.774 ± 0.107	0.772 ± 0.096	0.775 ± 0.101	0.781 ± 0.103
DT	G-mean	0.853 ± 0.081	0.854 ± 0.069	0.860 ± 0.067	0.864 ± 0.068
KNN	Accuracy	0.882 ± 0.085	0.883 ± 0.087	0.877 ± 0.087	0.881 ± 0.093
KNN	F-score	0.848 ± 0.116	0.849 ± 0.115	0.847 ± 0.118	0.852 ± 0.121
KNN	G-mean	0.896 ± 0.094	0.899 ± 0.090	0.904 ± 0.078	0.907 ± 0.080
LR	Accuracy	0.855 ± 0.074	0.870 ± 0.073	0.858 ± 0.077	0.870 ± 0.076
LR	F-score	0.812 ± 0.113	0.835 ± 0.105	0.825 ± 0.106	0.838 ± 0.106
LR	G-mean	0.875 ± 0.099	0.895 ± 0.075	0.899 ± 0.059	0.907 ± 0.059
RF	Accuracy	0.897 ± 0.080	0.905 ± 0.078	0.904 ± 0.078	0.906 ± 0.077
RF	F-score	0.867 ± 0.107	0.877 ± 0.103	0.875 ± 0.108	0.877 ± 0.108
RF	G-mean	0.911 ± 0.081	0.917 ± 0.078	0.923 ± 0.067	0.925 ± 0.065

Table 6: Optimal classification scores. The Maximum Performance (MP) classification scores are calculated using classifiers trained using the entire training set.

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

Classifier	Evaluation Metric	p-value	Significance
DT	Accuracy	2.1e-17	True
DT	F-score	2.5e-24	True
DT	G-mean	2.8e-16	True
KNN	Accuracy	1.1e-46	True
KNN	F-score	1.8e-66	True
KNN	G-mean	6.4e-42	True
LR	Accuracy	9.9e-59	True
LR	F-score	2.0e-76	True
LR	G-mean	2.2e-59	True
RF	Accuracy	5.7e-42	True
RF	F-score	4.6e-55	True
RF	G-mean	1.3e-38	True

Table 7: Results for Friedman test. 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.

618 Table 8 contains the *p-values* obtained with the Wilcoxon signed-rank
619 test. The proposed method was able to outperform both the standard AL
620 framework, as well as the AL framework using a normal oversampling policy
621 proposed in [19] with statistical significance in 9 out of 10 datasets.

622 The *p-values* shown in Table 9 refer to the results of the Holm-Bonferroni
623 test. The proposed method’s superior performance was statistically signifi-
624 cant for any combination of classifier and evaluation metric. Simultaneously,
625 the proposed method established statistical significance in the 3 scenarios
626 where the oversampling AL method failed to do so.

627 6.3. Discussion

628
629 In this paper we study the application of data augmentation methods
630 through the modification of the standard AL framework. This is done to
631 further reduce the amount of labeled data required to produce a reliable
632 classifier, at the expense of artificial data generation.

633 In Table 3 we found that the proposed method was able to outperform
634 the Standard AL framework in all scenarios. The mean rankings are consis-
635 tent with the mean AULC scores found in Table 4, while showing significant
636 performance differences between the proposed method and both the standard

Dataset	Oversampling	Standard
Baseball	5.0e-01	3.4e-01
Gas Drift	3.7e-26	4.6e-57
Image Segmentation	9.6e-18	2.1e-44
Japanese Vowels	2.4e-09	1.6e-32
Mfeat Zernike	1.2e-12	9.5e-40
Mice Protein	6.5e-32	1.5e-61
Pendigits	5.0e-18	2.3e-45
Texture	1.5e-22	6.7e-57
Vehicle	7.4e-11	7.9e-13
Waveform	8.9e-08	2.6e-02

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.

Classifier	Evaluation Metric	Oversampling	Proposed
DT	Accuracy	4.5e-05	1.6e-10
DT	F-score	1.9e-07	2.7e-10
DT	G-mean	2.5e-06	3.1e-09
KNN	Accuracy	5.5e-02	1.1e-05
KNN	F-score	6.7e-11	6.3e-14
KNN	G-mean	8.3e-06	1.3e-07
LR	Accuracy	8.1e-02	3.4e-06
LR	F-score	7.1e-06	2.0e-20
LR	G-mean	2.2e-07	1.1e-11
RF	Accuracy	2.0e-01	2.8e-02
RF	F-score	2.2e-05	8.1e-07
RF	G-mean	2.0e-04	2.0e-04

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).

637 and oversampling methods. The Friedman test in Table 7 showed that the
638 difference in the performance of these AL frameworks is statistically signifi-
639 cant, regardless of the classifier or performance metric being used.

640 The proposed method showed more consistent data utilization require-
641 ments to most of the assessed G-mean score thresholds when compared to
642 the remaining AL methods, as seen in Table 5. For example, to reach a
643 G-mean Score of 0.9 using the KNN and LR classifiers, the average amount
644 of data required with the Oversampling AL approach increased when com-
645 pared to the Standard approach. However, the proposed method was able to
646 decrease the amount of data required in both situations. The robustness of
647 the Proposed method is clearer in Figure 9. In most cases, this method was
648 able outperform the Oversampling method. At the same time, the proposed
649 method also addresses inconsistencies in situations where the Oversampling
650 method was unable to outperform the standard method.

651 The statistical analyses found in Tables 8 and 9 showed that the pro-
652 posed method’s superiority was statistically significant in all datasets except
653 one (Baseball) and established statistical significance when compared to the
654 Standard AL method for all combinations of classifier and performance met-
655 ric, including when the Oversampling AL method failed to do so. These
656 results show that the Proposed method increased the reliability of the new
657 AL framework and improved the quality of the final classifier while using less
658 data.

659 Even though it was not the core purpose of this study, we found that
660 the method proposed AL approach consistently outperformed the maximum
661 performance threshold. Specifically, in Table 6, the performance of the classi-
662 fiers originating from the proposed method was able to outperform classifiers
663 trained using the full training dataset in all 12 scenarios except one. This
664 suggests that the selection of a meaningful training subset training dataset
665 paired with data augmentation not only matches the classification perfor-
666 mance of ML algorithms, as it also improves them. Even in a setting with
667 fully labeled training data, the proposed method may be used as preprocess-
668 ing method to further optimize classification performance.

669 This study discussed the effect of data augmentation within the AL frame-
670 work, along with the exploration of optimal augmentation methods within
671 AL iterations. However, the conceptual nature of this study implies some
672 limitations. Specifically, the large amount of experiments required to test
673 the method’s efficacy, along with the limited computational power available,
674 led to a limited exploration of the grid search’s potential. Future work should

focus into 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 in 100% of scenarios. The exploration of other, more complex, data augmentation techniques might further improve its performance through the production of 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 into 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 found significant standard errors throughout our experimental results (see Subsection 6.1), which is consistent with the findings in [19, 24]. This 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.

7. Conclusion

The ability of training 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 to select the most informative observations and reduce the amount of required labeled data. On the other hand, the introduction of data variability in the training dataset can also be done via data augmentation. However, most, if not all, AL configurations using some form data augmentation are domain and/or task specific. These methods typically explore deep learning approaches on both classification and data augmentation. Consequently, they may not be applicable for 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 AL configuration, data selection in AL iterations aim 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 allows the reduction of some of the user interaction necessary to reach a sufficiently informative dataset. In order to further improve the proposed method future work will (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) better understand the gap between heuristic/input space data augmentation techniques and neural network/feature space data augmentation techniques in an AL context.

Declarations

Funding

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.

Code availability

The analyses and source code is available at github.com/joaopfonseca/ml-research.

References

- [1] V. Nath, D. Yang, B. A. Landman, D. Xu, H. R. Roth, Diminishing uncertainty within the training pool: Active learning for medical image segmentation, *IEEE Transactions on Medical Imaging* 40 (10) (2021) 2534–2547.
- [2] Y. Sverchkov, M. Craven, A review of active learning approaches to experimental design for uncovering biological networks, *PLoS Computational Biology* 13 (2017) e1005466.

- 744 [3] X. Li, D. Kuang, C. X. Ling, Active learning for hierarchical text classification,
745 Lecture Notes in Computer Science (including subseries Lecture
746 Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 7301
747 LNAI (2012) 14–25.
- 748 [4] Y. Li, J. Yin, L. Chen, Seal: Semisupervised adversarial active learning
749 on attributed graphs, IEEE Transactions on Neural Networks and
750 Learning Systems 32 (7) (2021) 3136–3147.
- 751 [5] O. Siméoni, M. Budnik, Y. Avrithis, G. Gravier, Rethinking deep active
752 learning: Using unlabeled data at model training, Proceedings -
753 International Conference on Pattern Recognition (2020) 1220–1227.
- 754 [6] J. E. Van Engelen, H. H. Hoos, A survey on semi-supervised learning,
755 Machine Learning 109 (2) (2020) 373–440.
- 756 [7] O. Sener, S. Savarese, Active learning for convolutional neural networks:
757 A core-set approach, in: International Conference on Learning Representations,
758 2018.
- 759 [8] Y. Leng, X. Xu, G. Qi, Combining active learning and semi-supervised
760 learning to construct svm classifier, Knowledge-Based Systems 44 (2013)
761 121–131.
- 762 [9] H. Yu, X. Yang, S. Zheng, C. Sun, Active learning from imbalanced
763 data: A solution of online weighted extreme learning machine, IEEE
764 Transactions on Neural Networks and Learning Systems 30 (2019) 1088–
765 1103.
- 766 [10] W. Zong, G.-B. Huang, Y. Chen, Weighted extreme learning machine
767 for imbalance learning, Neurocomputing 101 (2013) 229–242.
- 768 [11] J. Qin, C. Wang, Q. Zou, Y. Sun, B. Chen, Active learning with extreme
769 learning machine for online imbalanced multiclass classification,
770 Knowledge-Based Systems 231 (2021) 107385.
- 771 [12] W. Liu, H. Zhang, Z. Ding, Q. Liu, C. Zhu, A comprehensive active
772 learning method for multiclass imbalanced data streams with concept
773 drift, Knowledge-Based Systems 215 (2021) 106778.

- 774 [13] A. Tharwat, W. Schenck, Balancing exploration and exploitation: A
775 novel active learner for imbalanced data, *Knowledge-Based Systems* 210
776 (2020) 106500.
- 777 [14] Y.-Y. Kim, K. Song, J. Jang, I.-c. Moon, Lada: Look-ahead data ac-
778 quisition via augmentation for deep active learning, *Advances in Neural*
779 *Information Processing Systems* 34 (2021).
- 780 [15] G. Douzas, F. Bacao, Geometric SMOTE a geometrically enhanced
781 drop-in replacement for SMOTE, *Information Sciences* 501 (2019) 118–
782 135.
- 783 [16] P. Ren, Y. Xiao, X. Chang, P.-Y. Huang, Z. Li, B. B. Gupta, X. Chen,
784 X. Wang, A survey of deep active learning, *ACM Computing Surveys*
785 (CSUR) 54 (9) (2021) 1–40.
- 786 [17] V. K. Shrivastava, M. K. Pradhan, Hyperspectral remote sensing im-
787 age classification using active learning, *Studies in Computational Intel-*
788 *ligence* 907 (2021) 133–152.
- 789 [18] P. Ren, Y. Xiao, X. Chang, P.-Y. Huang, Z. Li, X. Chen, X. Wang, A
790 survey of deep active learning, *arXiv preprint arXiv:2009.00236* (2020).
- 791 [19] J. Fonseca, G. Douzas, F. Bacao, Increasing the Effectiveness of Active
792 Learning: Introducing Artificial Data Generation in Active Learning
793 for Land Use/Land Cover Classification, *Remote Sensing 2021*, Vol. 13,
794 Page 2619 13 (13) (2021) 2619.
- 795 [20] D. Yoo, I. S. Kweon, Learning loss for active learning, in: *Proceedings*
796 *of the IEEE/CVF Conference on Computer Vision and Pattern Recog-*
797 *nition*, 2019, pp. 93–102.
- 798 [21] H. H. Aghdam, A. Gonzalez-Garcia, A. Lopez, J. Weijer, Active learn-
799 ing for deep detection neural networks, in: *Proceedings of the IEEE*
800 *International Conference on Computer Vision*, Vol. 2019-Octob, 2019,
801 pp. 3671–3679.
- 802 [22] T. Su, S. Zhang, T. Liu, Multi-spectral image classification based on an
803 object-based active learning approach, *Remote Sensing* 12 (2020) 504.

- [23] Z. del Rosario, M. Rupp, Y. Kim, E. Antono, J. Ling, Assessing the frontier: Active learning, model accuracy, and multi-objective candidate discovery and optimization, *The Journal of Chemical Physics* 153 (2020) 024112.
- [24] 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.
- [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 (2020) 1–17.
- [26] H. T. Nguyen, A. Smeulders, Active learning using pre-clustering, in: *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 79.
- [27] B. Gu, Z. Zhai, C. Deng, H. Huang, Efficient active learning by querying discriminative and representative samples and fully exploiting unlabeled data, *IEEE Transactions on Neural Networks and Learning Systems* 32 (2021) 4111–4122.
- [28] P. Kumar, A. Gupta, Active learning query strategies for classification, regression, and clustering: A survey, *Journal of Computer Science and Technology* 2020 35:4 35 (2020) 913–945.
- [29] Y. Fu, X. Zhu, B. Li, A survey on instance selection for active learning, *Knowledge and information systems* 35 (2) (2013) 249–283.
- [30] A. Samat, P. Gamba, S. Liu, P. Du, J. Abuduwaili, Jointly informative and manifold structure representative sampling based active learning for remote sensing image classification, *IEEE Transactions on Geoscience and Remote Sensing* 54 (2016) 6803–6817.
- [31] S. J. Huang, R. Jin, Z. H. Zhou, Active learning by querying informative and representative examples, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36 (2014) 1936–1949.
- [32] D. Ienco, A. Bifet, I. Žliobaite, B. Pfahringer, Clustering based active learning for evolving data streams, *Lecture Notes in Computer Science*

- (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 8140 LNAI (2013) 79–93.
- [33] K. Brinker, Incorporating diversity in active learning with support vector machines, in: Proceedings of the 20th international conference on machine learning (ICML-03), 2003, pp. 59–66.
- [34] J. Jia, M. T. Schaub, S. Segarra, A. R. Benson, Graph-based semi-supervised & active learning for edge flows, in: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2019, pp. 761–771.
- [35] B. Settles, From theories to queries: Active learning in practice, in: Active Learning and Experimental Design workshop In conjunction with AISTATS 2010, JMLR Workshop and Conference Proceedings, 2011, pp. 1–18.
- [36] D. D. Lewis, W. A. Gale, A sequential algorithm for training text classifiers, in: SIGIR’94, Springer, 1994, pp. 3–12.
- [37] T. Luo, K. Kramer, D. B. Goldgof, L. O. Hall, S. Samson, A. Remsen, T. Hopkins, D. Cohn, Active learning to recognize multiple types of plankton., Journal of Machine Learning Research 6 (4) (2005).
- [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) 1092.
- [39] 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) (2012) 844–856.
- [40] N. Abe, Query learning strategies using boosting and bagging, Proc. of 15th Int. Conf. on Machine Learning (ICML98) (1998) 1–9.
- [41] P. Melville, R. J. Mooney, Diverse ensembles for active learning, in: Proceedings of the twenty-first international conference on Machine learning, 2004, p. 74.

- [42] M. Bloodgood, Support vector machine active learning algorithms with query-by-committee versus closest-to-hyperplane selection, Proceedings - 12th IEEE International Conference on Semantic Computing, ICSC 2018 2018-January (2018) 148–155.
- [43] J. Zhou, S. Sun, Improved margin sampling for active learning, Communications in Computer and Information Science 483 (2014) 120–129.
- [44] S. Behpour, K. M. Kitani, B. D. Ziebart, Ada: Adversarial data augmentation for object detection, Proceedings - 2019 IEEE Winter Conference on Applications of Computer Vision, WACV 2019 (2019) 1243–1252.
- [45] Z. Zhong, L. Zheng, G. Kang, S. Li, Y. Yang, Random erasing data augmentation, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 13001–13008.
- [46] 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.
- [47] C. Shorten, T. M. Khoshgoftaar, A survey on image data augmentation for deep learning, Journal of Big Data 6 (1) (2019) 1–48.
- [48] B. K. Iwana, S. Uchida, An empirical survey of data augmentation for time series classification with neural networks, Plos one 16 (7) (2021) e0254841.
- [49] S. C. Wong, A. Gatt, V. Stamatescu, M. D. McDonnell, Understanding data augmentation for classification: when to warp?, in: 2016 international conference on digital image computing: techniques and applications (DICTA), IEEE, 2016, pp. 1–6.
- [50] O. Kashefi, 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), Association for Computational Linguistics, Online, 2020, pp. 200–208.
- [51] 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.

- 898 [52] J. Fonseca, G. Douzas, F. Bacao, Improving imbalanced land cover clas-
899 sification with k-means smote: Detecting and oversampling distinctive
900 minority spectral signatures, *Information* 12 (7) (2021) 266.
- 901 [53] G. Douzas, F. Bacao, J. Fonseca, M. Khudinyan, Imbalanced learning
902 in land cover classification: Improving minority classes' prediction ac-
903 curacy using the geometric smote algorithm, *Remote Sensing* 11 (24)
904 (2019) 3040.
- 905 [54] H. Han, W.-Y. Wang, B.-H. Mao, Borderline-smote: A new over-
906 sampling method in imbalanced data sets learning, in: *International*
907 *Conference on Intelligent Computing*, Springer, Berlin, Heidelberg,
908 2005, pp. 878–887.
- 909 [55] G. Douzas, F. Bacao, F. Last, Improving imbalanced learning through
910 a heuristic oversampling method based on k-means and smote, *Informa-*
911 *tion Sciences* 465 (2018) 1–20.
- 912 [56] Y. Ma, S. Lu, E. Xu, T. Yu, L. Zhou, Combining active learning and
913 data augmentation for image classification, in: *Proceedings of the 2020*
914 *3rd International Conference on Big Data Technologies*, 2020, pp. 58–62.
- 915 [57] H. Quteineh, S. Samothrakis, R. Sutcliffe, Textual data augmentation
916 for efficient active learning on tiny datasets, in: *Proceedings of the*
917 *2020 Conference on Empirical Methods in Natural Language Processing*
918 *(EMNLP)*, 2020, pp. 7400–7410.
- 919 [58] Q. Li, Z. Huang, Y. Dou, Z. Zhang, A framework of data augmentation
920 while active learning for chinese named entity recognition, in: *Interna-*
921 *tional Conference on Knowledge Science, Engineering and Management*,
922 Springer, 2021, pp. 88–100.
- 923 [59] C. Wu, *The decision tree approach to classification.*, Purdue University,
924 1975.
- 925 [60] T. Cover, P. Hart, Nearest neighbor pattern classification, *IEEE Trans-*
926 *actions on Information Theory* 13 (1) (1967) 21–27.
- 927 [61] T. K. Ho, Random decision forests, in: *Proceedings of the Third Inter-*
928 *national Conference on Document Analysis and Recognition (Volume 1)*
929 *- Volume 1, ICDAR '95, IEEE Computer Society, USA, 1995*, p. 278.

- 930 [62] J. A. Nelder, R. W. Wedderburn, Generalized linear models, *Journal of*
931 *the Royal Statistical Society: Series A (General)* 135 (3) (1972) 370–384.
- 932 [63] L. A. Jeni, J. F. Cohn, F. De La Torre, Facing imbalanced data - Recom-
933 mendations for the use of performance metrics, in: *Proceedings - 2013*
934 *Humaine Association Conference on Affective Computing and Intelligent*
935 *Interaction, ACII 2013, 2013*, pp. 245–251.
- 936 [64] M. Fatourechi, R. K. Ward, S. G. Mason, J. Huggins, A. Schloegl, G. E.
937 Birch, Comparison of evaluation metrics in classification applications
938 with imbalanced datasets, in: *2008 seventh international conference on*
939 *machine learning and applications, IEEE, 2008*, pp. 777–782.
- 940 [65] M. Kubat, S. Matwin, et al., Addressing the curse of imbalanced training
941 sets: one-sided selection, in: *Icml, Vol. 97, Citeseer, 1997*, pp. 179–186.
- 942 [66] T. Reitmaier, B. Sick, Let us know your decision: Pool-based active
943 training of a generative classifier with the selection strategy 4ds, *Infor-*
944 *mation Sciences* 230 (2013) 106–131.
- 945 [67] J.-F. Kagy, T. Kayadelen, J. Ma, A. Rostamizadeh, J. Strnadova, The
946 practical challenges of active learning: Lessons learned from live exper-
947 imentation, *arXiv preprint arXiv:1907.00038* (6 2019).
- 948 [68] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion,
949 O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vander-
950 plas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, É. Duchesnay,
951 *Scikit-learn: Machine Learning in Python, Journal of Machine Learning*
952 *Research* 12 (Oct) (2011) 2825–2830.
- 953 [69] G. Lemaître, F. Nogueira, C. K. Aridas, Imbalanced-learn: A python
954 toolbox to tackle the curse of imbalanced datasets in machine learning,
955 *Journal of Machine Learning Research* 18 (17) (2017) 1–5.
- 956 [70] J. Demšar, Statistical comparisons of classifiers over multiple data sets,
957 *Journal of Machine Learning Research* 7 (2006) 1–30.
- 958 [71] M. Friedman, The use of ranks to avoid the assumption of normality
959 implicit in the analysis of variance, *Journal of the american statistical*
960 *association* 32 (200) (1937) 675–701.

- 961 [72] F. Wilcoxon, Individual Comparisons by Ranking Methods, Biometrics
962 Bulletin 1 (6) (1945) 80.
- 963 [73] S. Holm, A simple sequentially rejective multiple test procedure, Scan-
964 dinavian journal of statistics (1979) 65–70.