



Imbalanced Data Stream Classification Using Hybrid Data Preprocessing

Barbara Bobowska^(✉), Jakub Klikowski^(✉), and Michał Woźniak^(✉)

Department of Systems and Computer Networks,
Wrocław University of Science and Technology, Wrocław, Poland
{barbara.bobowska, jakub.klikowski, michal.wozniak}@pwr.edu.pl

Abstract. Imbalanced data streams have gained significant popularity among the researchers in recent years. This area of research is not only still greatly underdeveloped, but there are also numerous inherent difficulties that need to be addressed when creating algorithms that could be utilized in such dynamic environment and achieve satisfactory results when it comes to their predictive abilities. In this paper, a novel algorithm that combines both over- and under-sampling techniques in order to create a more robust classifier dedicated to imbalanced data streams is proposed. The efficiency and high predictive quality of the proposed method have been confirmed on the basis of extensive experimental research carried out on the real and the computer-generated data streams.

Keywords: Imbalanced data · Data stream classification · Data preprocessing

1 Introduction

In the last couple of years, a sharp rise in products and systems using machine learning to enhance their performance is observed. Many of the applications such as predicting user behavior on social platforms like Twitter, or client activity on online stores fall into the category of imbalanced data stream classification [24]. When designing methods for data stream classification one has to take into account the characteristics of a data stream such as the sequential manner that the data arrives, over which one has no control when it comes to the order of the arriving samples, as well as the fact that the size of the stream could be possibly infinite. Due to that requirement, it is impossible to process the upcoming data in multiple passes and such the samples can be processed once [26]. Furthermore, one has to consider the rapid rate at which the data arrives, at the same time ensuring that the processing of the data stream is done in a timely fashion such that the delay in the performance of the algorithm is minimal. Data streams can exhibit a change in data and target concepts over time (so-called Non-stationary data streams) [16, 26]. Such a phenomenon is called *concept drift* [12] and it is quite common i.e. the change of popular topics on Twitter.

Due to the *concept drift* the performance of the classifier can degrade over time and as such the classifier has to be trained incrementally to accommodate the changes of concepts of non-stationary data streams. Moreover, the proportion between classes is often skewed with one class being over-represented. In cases where the imbalance ratio is present traditional accuracy driven methods are not applicable especially when misclassification of the minority class examples is much more costly, as is often the case i.e. fraud detection [24]. It is worth mentioning that not only the imbalance ratio can influence the performance of the classifier. Some examples can be easy to classify even when the IR is high if the classes are separated from each other the decision boundary and be determined with ease. However, it has been observed that instances of the minority class have a tendency to create sparsely spread throughout the object space clusters, often surrounded by majority class examples [4]. The presence of noise and outliers is another difficulty factor that needs to be addressed. In [3,15] authors created preprocessing methods with those issues in mind.

Data streams may be processed either in blocks or one instance at a time. One of the most important issues in learning from the data stream is when to update the classifier [22]. Most researchers distinguish between two approaches: active and passive. In the former, the update is performed only if drift is detected while the later updates the classifier continuously regardless if the drift was detected or not [9]. In order to satisfy the time and memory requirements, a forgetting or data management mechanism must be used. One of the most popular approaches to forgetting is using sliding windows, which can be either sequence based, where the size of the window is defined by a number of instances and time stamp based where the window is defined by a certain duration time. In the simplest example sliding windows are of fixed size, and include only the most recent examples. The oldest samples in a window are discarded in favor of new ones. Some methods implement sliding windows of varying size depending on the response from drift detectors [2].

The main contributions of this work are as follows:

- Proposition of the two novel imbalanced data stream classifiers (DSC-R and DSC-S) which employ under- and oversampling techniques for balancing data.
- Experimental evaluation of the proposed algorithms and their comparison with state-of-art methods.

The article is organized as follows. Sections 1 and 2 present a brief introduction to the problem of imbalance data stream classification and a quick overview of the state-of-the-art algorithms dedicated to it. Section 3 offers an in-depth explanation of the proposed solution. Section 4 showcases the results of the computer experiments, comparing the proposed algorithm to different techniques for imbalanced data classification, proving the usefulness of the developed algorithm. Section 5 presents the conclusions and describes possible future improvements to the proposed method.

2 Related Works

Studies over the years presented algorithms dedicated to data stream analysis. Very fast decision tree (VFDT) proposed by Domingos and Hulten [13] was among the first methods for stream analysis, that to this day has been a basis for many modifications. VFDT utilizes the Hoeffding bound in order to calculate the proper number of examples needed to select the split-node. The algorithm incrementally creates a tree form from a data stream ensuring that once the examples were used to update the tree they are negligible and can be removed. The aforementioned modifications include ideas such as pruning mechanisms or utilizing sliding windows and drift detectors in order to better the algorithms in case of non-stationary streams [10]. Worth noting are several methods using ensembles of classifiers. Weighted Majority Algorithm [18] adjusts the weights of the classifiers in the ensemble so that the weight of an expert that misclassified an instance is decreased accordingly to the user-specified value. A modification of the method with an added procedure which adds new classifiers to the ensemble when the overall performance is unsatisfactory called Dynamic Weighted Majority (DWM) was introduced in [14]. In Accuracy Weighted Ensemble (AWE) a new classifier is added only if the ensemble's size is not exceeded [25] while in Learn++.NSE [8] such a constraint is not applied. In Learn++.CDS Ditzler and Polikar combine their previous work Learn++.NSE with SMOTE sampling in order to better address the data imbalance and later replacing SMOTE with an original bagging-based method of data balancing [7]. In SEA [23] a new classifier candidate is evaluated to determine whether or not it is worth including into the ensemble at the cost of replacing some other classifier already in the ensemble. Other approaches such as OUSEnsemble (Over Under Sampling Ensemble) [11] make use of sampling techniques. The stream is divided into blocks that consist of examples from both majority and minority class. The idea is to propagate all the instances of the minority class from the previous block and under-sample the majority examples in the current block such that the desired imbalance ratio is acquired. Afterwards, from the resultant subset, datasets later used to build component classifiers for the ensemble, are created by propagating all instances of the minority class to each of the datasets while each example from the majority class is propagated to only one dataset. Proposed by Chen and He the Selectively Recursive Approach (SERA) [5] uses a Mahalanobis distance to determine which of the examples from the minority class in the previous block are most similar to the minority examples in the current block. Based on that a limited number of minority class examples is selected and added to the majority class examples in the current block. Chen and He later designed a Recursive Ensemble Approach (REA) [6]. In REA minority class examples from the previous block that are nearest neighbors of minority class examples from the current block are added in order to balance the given training block. Both REA and SERA proved to make more accurate predictions than the method proposed by [19]. A Chunk-based ensemble approach, proposed by Wang et al. called KMean-Clustering [25] utilizes k-mean clustering in order to under-sample the majority class, by using the centroids created in the clustering process to resample the majority instances.

3 The Deterministic Sampling Classifier

The proposed method, called *Deterministic Sampling Classifier* (DSC), for data stream classification, processes the upcoming data in chunks. Each chunk is used in two operations. Firstly, the instances of the majority class present in the currently processed block are under-sampled in order to produce a balanced class representation in a data chunk (Fig. 1).

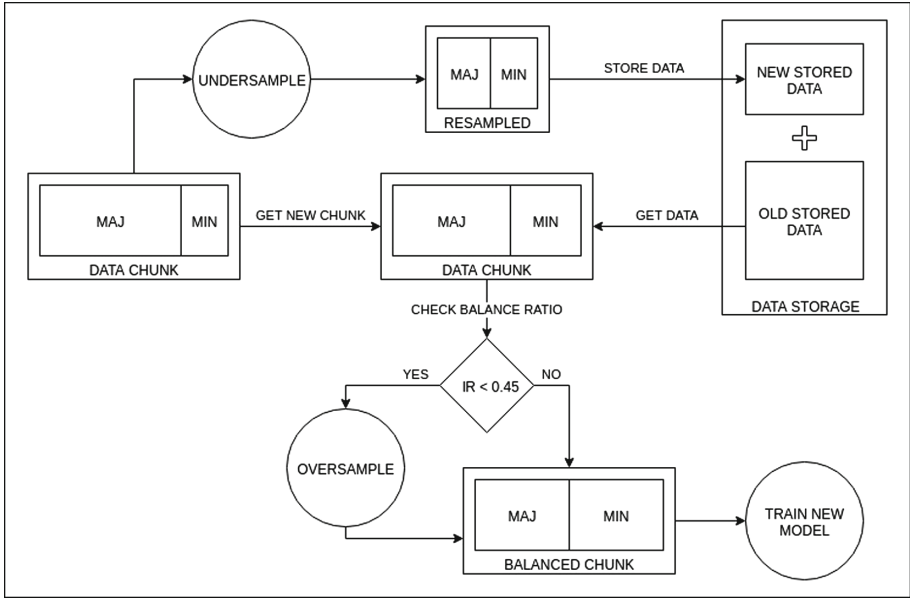


Fig. 1. Proposed method flow diagram

The resulting data (referred to in the Fig. 1 as NEW STORED DATA) is then stored in a memory buffer (DATA STORAGE). Secondly, that same block of data is combined with a part of the data from the buffer, called OLD STORED DATA, using GET NEW CHUNK, which copies the data from the currently processed block and the GET DATA method, which takes OLD STORED DATA from the DATA STORAGE buffer. OLD STORED DATA, consists of all the previously accumulated under-sampled chunks. When a new chunk of data arrives the data from NEW STORED DATA are moved to the OLD STORED DATA part of the buffer. The DATA STORAGE is of fixed size. When the buffer is full, the oldest examples are removed from it. Afterward, the imbalance ratio of the data block created as a result of the GET NEW CHUNK and GET DATA is calculated, and if the value is lower than 0.45, an oversampling of the minority class is performed, and then a classification model is trained. Otherwise, the algorithm accepts the chunk as properly balanced and uses it to train the

model right away. The implementation allows one to choose sampling algorithms of their liking. In this paper, the authors created two versions of the method DSC-S (Deterministic Sampling Classifier-SMOTE) and DSC-R (Deterministic Sampling Classifier-Random). For the DSC-R method the chosen sampling methods were: random over- and under-sampling and for the DSC-S: SMOTE and NearMiss (implementation from the imbalanced-learn library [17]) for over- and under-sampling accordingly.

4 Experimental Evaluation

The quality of the proposed algorithms was evaluated on the basis of computer experiments, using 26 real and 60 synthetic data streams. The evaluation procedure used in order to assess the predictive performance of a data stream classifier was conducted by interleaving testing with training (test-then-train) [16]. Each block is first used to test the classifier and afterward it is used for training. As a measure of comparison, the following methods were used: OUSEnsemble, KMeanClustering, REA, Learn++.CDS, Learn++.NIE and MLPClassifier (Multi-layer Perceptron classifier), using a k-NN as a base classifier. The algorithms were implemented in Python using Scikit-learn [20] and imbalanced-learn [17] libraries¹. The selected real streams were downloaded from the KEEL [1] and PROMISE Software Engineering Repository [21]. The chosen datasets consisted of multidimensional binary classification problems with the imbalance ratio ranging from 1 to 39. The datasets were described in Table 1. The results were analyzed using the KEEL software evaluation tool [1]. Non-parametrical statistical tests were performed namely the Friedman Test as well as a Nemenyi's Post-Hoc Procedure. The metrics chosen were F-score, Gmean and AUC score. Tables 2, 3 and 4 present the obtained results. The table presents the obtained results as the mean value for each of the metrics, as well as, the information on those methods that performed poorly in comparison with the method named in the column, placed directly below the score. For instance, given the abalone-17-vs-7-8-9-10 dataset, the DSC-R algorithm performed statistically better than the 3rd, 5th, the 7th and 8th algorithm in the table (read from left to right). The obtained results prove the usefulness of the proposed algorithms. For the F-score the proposed DSC-R and DSC-S algorithms along with the REA algorithm have the best results. What is interesting the MLPC algorithm performed consistently the worst. For the Gmean the results are similar. The methods introduced in the paper perform favorably in comparison with other algorithms, the Learn++.CDS and Learn++.NIE techniques, as well as REA, have comparable results to the DSC-R and DSC-S methods. Lastly, the results in Table 4 representing the results for AUC score indicate the proposed algorithm obtained satisfactory results,

¹ Repository link: <https://github.com/w4k2/iot-ecml2019>.

with only the LCDS algorithm performing marginally better. It is worth mentioning, that the created methods are robust enough, so that imbalance ratio (whether low or high) does not negatively impact their performance.

Table 1. Overview of datasets selected for experimental evaluation (source: KEEL and PROMISE Software Engineering Repository).

Dataset	IMB. RATIO	SAMPLES	FEATURES
<i>abalone-17_vs-7-8-9-10</i>	39	2338	8
<i>australian</i>	1	690	14
<i>glass-0-1-2-3-vs-4-5-6</i>	3	214	9
<i>glass0</i>	2	214	9
<i>glass1</i>	2	214	9
<i>heart</i>	1	270	13
<i>jm1</i>	4	10885	22
<i>kc1</i>	5	2109	22
<i>kc2</i>	4	522	22
<i>kr-vs-k-three_vs_eleven</i>	35	2935	6
<i>kr-vs-k-zero-one_vs_draw</i>	26	2901	6
<i>page-blocks0</i>	9	5472	10
<i>pima</i>	2	768	8
<i>segment0</i>	6	2308	19
<i>shuttle-c0-vs-c4</i>	14	1829	9
<i>vehicle0</i>	3	846	18
<i>vowel0</i>	10	988	13
<i>wisconsin</i>	2	683	9
<i>yeast-0-2-5-6-vs-3-7-8-9</i>	9	1004	8
<i>yeast-0-2-5-7-9-vs-3-6-8</i>	9	1004	8
<i>yeast-0-3-5-9-vs-7-8</i>	9	506	8
<i>yeast-0-5-6-7-9-vs-4</i>	9	528	8
<i>yeast-2-vs-4</i>	9	514	8
<i>yeast1</i>	2	1484	8
<i>yeast3</i>	8	1484	8

Table 2. Overview of the results for F-score.

Dataset	1	2	3	4	5	6	7	8
	DSC-R	DSC-S	KMC	LCDS	LNIE	REA	OUSE	MLPC
<i>abalone-17-vs-7-8-9-10</i>	0.237 3,5,7,8	0.208 3,5,7,8	0.060 8	0.250 3,5,7,8	0.046 –	0.176 –	0.072 8	0.000 –
<i>australian</i>	0.571 3,8	0.601 3,6,8	0.444 8	0.514 8	0.604 3,4,6,8	0.518 8	0.602 3,4,6,8	0.000 –
<i>electricity-normalized</i>	0.631 3,4,5,6,7,8	0.623 3,4,6,7,8	0.580 8	0.588 8	0.598 8	0.575 8	0.593 8	0.324 –
<i>glass-0-1-2-3-vs-4-5-6</i>	0.765 7,8	0.801 7,8	0.781 7,8	0.814 7,8	0.768 7,8	0.782 7,8	0.452 8	0.000 –
<i>glass0</i>	0.662 5,7,8	0.673 5,7,8	0.592 7,8	0.640 7,8	0.516 8	0.597 8	0.473 8	0.000 –
<i>glass1</i>	0.687 4,5,6,8	0.649 6,8	0.576 6,8	0.517 8	0.558 6,8	0.412 8	0.573 6,8	0.000 –
<i>heart</i>	0.583 8	0.569 8	0.489 8	0.569 8	0.641 3,8	0.538 8	0.580 8	0.144 –
<i>jm1</i>	0.395 3,7,8	0.390 3,7,8	0.182 –	0.399 3,6,7,8	0.372 3,8	0.331 3,8	0.317 3,8	0.170 –
<i>kc1</i>	0.387 3,7,8	0.394 3,6,7,8	0.186 –	0.377 3,7,8	0.394 3,7,8	0.331 3,8	0.311 3,8	0.127 –
<i>kc2</i>	0.532 3	0.560 3,7,8	0.198 –	0.511 3	0.509 3	0.467 3	0.380 3	0.389 3
<i>kr-vs-k-three-vs-eleven</i>	0.701 3,5,7,8	0.847 1,3,4,5,7,8	0.256 8	0.730 3,5,7,8	0.408 3,8	0.836 1,3,4,5,7,8	0.335 8	0.000 –
<i>kr-vs-k-zero-one-vs-draw</i>	0.673 3,7,8	0.804 1,3,5,7,8	0.395 8	0.773 1,3,5,7,8	0.591 3,7,8	0.785 1,3,5,7,8	0.367 8	0.000 –
<i>page-blocks0</i>	0.541 3	0.576 3,4,5,7,8	0.252 –	0.359 –	0.335 –	0.533 3,7	0.352 –	0.323 –
<i>pima</i>	0.601 3,5,7,8	0.578 8	0.508 8	0.539 8	0.536 8	0.581 8	0.518 8	0.370 –
<i>segment0</i>	0.671 3,4,5,6,7,8	0.792 1,3,4,5,6,7,8	0.405 7,8	0.587 3,5,7,8	0.339 8	0.580 3,5,7,8	0.294 8	0.123 –
<i>shuttle-c0-vs-c4</i>	0.995 3,5,7,8	0.995 3,5,7,8	0.923 8	0.960 8	0.931 8	0.995 3,5,7,8	0.955 8	0.202 –
<i>vehicle0</i>	0.812 3,5,7,8	0.824 3,4,5,7,8	0.653 8	0.744 5,7,8	0.563 8	0.758 5,7,8	0.563 8	0.257 –
<i>vowel0</i>	0.682 3,5,7,8	0.753 3,4,5,6,7,8	0.274 8	0.553 3,5,7,8	0.272 8	0.553 3,5,7,8	0.251 8	0.059 –
<i>wisconsin</i>	0.966 3,8	0.951 3,8	0.708 8	0.942 3,8	0.941 3,8	0.949 3,8	0.862 3,8	0.019 –
<i>yeast-0-2-5-6-vs-3-7-8-9</i>	0.446 7,8	0.393 8	0.358 8	0.362 8	0.389 8	0.514 3,4,7,8	0.257 8	0.073 –
<i>yeast-0-2-5-7-9-vs-3-6-8</i>	0.667 7,8	0.710 3,5,7,8	0.580 7,8	0.675 7,8	0.590 7,8	0.722 3,5,7,8	0.366 8	0.000 –
<i>yeast-0-3-5-9-vs-7-8</i>	0.251 8	0.288 8	0.222 8	0.265 8	0.145 –	0.403 3,5,7,8	0.185 8	0.018 –
<i>yeast-0-5-6-7-9-vs-4</i>	0.350 7,8	0.432 7,8	0.327 7,8	0.428 7,8	0.257 8	0.379 7,8	0.176 8	0.000 –
<i>yeast-2-vs-4</i>	0.615 7,8	0.670 7,8	0.581 8	0.569 8	0.571 8	0.702 7,8	0.389 8	0.000 –
<i>yeast1</i>	0.527 8	0.495 8	0.520 8	0.504 8	0.430 8	0.476 8	0.498 8	0.000 –
<i>yeast3</i>	0.566 7,8	0.589 7,8	0.429 8	0.573 7,8	0.436 8	0.644 7,8	0.314 8	0.100 –

Table 3. Overview of the results for Gmean score.

Dataset	1	2	3	4	5	6	7	8
	DSC-R	DSC-S	KMC	LCDS	LNIE	REA	OUSE	MLPC
<i>abalone-17-vs-7-8-9-10</i>	0.782 3,5,6,7,8	0.818 3,5,6,7,8	0.541 5,8	0.684 5,6,8	0.102 -	0.282 -	0.536 5,8	0.000 -
<i>australian</i>	0.622 3,5,7,8	0.640 3,4,5,6,7,8	0.433 5,7,8	0.557 3,5,7,8	0.187 8	0.580 3,5,7,8	0.100 -	0.000 -
<i>electricity-normalized</i>	0.682 3,4,5,6,7,8	0.676 3,4,5,6,7,8	0.639 5,7,8	0.644 5,7,8	0.510 7,8	0.635 5,7,8	0.284 -	0.393 7
<i>glass-0-1-2-3-vs-4-5-6</i>	0.865 7,8	0.892 7,8	0.836 7,8	0.905 7,8	0.869 7,8	0.869 7,8	0.165 -	0.000 -
<i>glass0</i>	0.739 3,5,7,8	0.756 3,5,7,8	0.588 5,7,8	0.720 5,7,8	0.287 8	0.697 5,7,8	0.071 -	0.000 -
<i>glass1</i>	0.750 3,4,5,6,7,8	0.717 3,4,5,6,7,8	0.560 5,7,8	0.600 5,7,8	0.254 8	0.526 5,7,8	0.127 -	0.000 -
<i>heart</i>	0.623 3,8	0.618 3,8	0.518 8	0.589 8	0.346 -	0.566 8	0.284 -	0.205 -
<i>jm1</i>	0.660 3,6,7,8	0.666 3,6,7,8	0.275 -	0.675 3,6,7,8	0.630 3,7,8	0.557 3,7,8	0.369 -	0.291 -
<i>kc1</i>	0.651 3,6,7,8	0.670 3,6,7,8	0.278 -	0.655 3,6,7,8	0.658 3,6,7,8	0.557 3,7,8	0.391 8	0.204 -
<i>kc2</i>	0.717 3,7,8	0.772 3,7,8	0.230 -	0.711 3,7,8	0.715 3,7,8	0.618 3,7	0.259 -	0.347 -
<i>kr-vs-k-three-vs-eleven</i>	0.988 3,7,8	0.985 3,7,8	0.911 8	0.981 3,7,8	0.836 8	0.985 3,7,8	0.933 8	0.000 -
<i>kr-vs-k-zero-one-vs-draw</i>	0.970 3,5,7,8	0.965 3,5,8	0.898 8	0.932 8	0.873 8	0.933 8	0.935 8	0.000 -
<i>page-blocks0</i>	0.850 3,4,5,8	0.848 3,4,5,8	0.543 -	0.586 -	0.649 -	0.812 3,8	0.805 3,8	0.447 -
<i>pima</i>	0.690 3,4,5,7,8	0.673 3,5,7,8	0.590 5,7,8	0.634 5,7,8	0.400 7	0.670 3,5,7,8	0.197 -	0.482 7
<i>segment0</i>	0.911 3,4,5,6,7,8	0.934 3,4,5,6,7,8	0.719 5,7,8	0.848 3,5,7,8	0.585 7,8	0.851 3,5,7,8	0.415 8	0.242 -
<i>shuttle-c0-vs-c4</i>	0.995 8	0.995 8	0.949 8	0.963 8	0.950 8	0.995 8	0.996 3,5,8	0.304 -
<i>vehicle0</i>	0.916 3,5,7,8	0.901 3,5,7,8	0.819 5,7,8	0.881 5,7,8	0.723 8	0.882 5,7,8	0.708 -	0.380 -
<i>vowel0</i>	0.937 3,4,5,7,8	0.938 3,4,5,7,8	0.523 8	0.797 3,5,7,8	0.634 8	0.790 3,8	0.634 8	0.224 -
<i>wisconsin</i>	0.974 3,8	0.960 3,8	0.747 8	0.953 3,8	0.950 3,8	0.957 3,8	0.781 8	0.042 -
<i>yeast-0-2-5-6-vs-3-7-8-9</i>	0.761 7,8	0.661 8	0.685 7,8	0.709 7,8	0.683 8	0.745 7,8	0.511 8	0.178 -
<i>yeast-0-2-5-7-9-vs-3-6-8</i>	0.862 7,8	0.848 7,8	0.827 7,8	0.879 7,8	0.839 7,8	0.875 7,8	0.681 8	0.000 -
<i>yeast-0-3-5-9-vs-7-8</i>	0.617 8	0.650 7,8	0.517 8	0.596 8	0.339 -	0.651 8	0.470 8	0.069 -
<i>yeast-0-5-6-7-9-vs-4</i>	0.750 7,8	0.822 7,8	0.777 7,8	0.783 7,8	0.612 8	0.694 8	0.485 8	0.000 -
<i>yeast-2-vs-4</i>	0.827 8	0.826 8	0.801 8	0.830 8	0.723 8	0.847 8	0.731 8	0.000 -
<i>yeast1</i>	0.634 5,7,8	0.616 5,7,8	0.599 5,7,8	0.612 5,7,8	0.330 8	0.593 5,7,8	0.344 8	0.000 -
<i>yeast3</i>	0.871 7,8	0.855 7,8	0.815 8	0.841 7,8	0.773 8	0.875 5,7,8	0.698 8	0.197 -

Table 4. Overview of the results for AUC score.

Dataset	1	2	3	4	5	6	7	8
	DSC-R	DSC-S	KMC	LCDS	LNIE	REA	OUSE	MLPC
<i>abalone-17-vs-7-8-9-10</i>	0.997 3,5,6,7,8	0.989 3,5,6,7,8	0.684 8	0.999 3,5,6,7,8	0.828 3,8	0.967 3,5,7,8	0.789 3,8	0.329 –
<i>australian</i>	0.827 3,5,7,8	0.829 3,5,7,8	0.551 8	0.819 3,5,7,8	0.662 3,7,8	0.803 3,5,7,8	0.578 8	0.448 –
<i>electricity-normalized</i>	0.961 3,5,7,8	0.961 3,5,7,8	0.927 5,7,8	0.974 1,2,3,5,6,7,8	0.871 7,8	0.971 1,2,3,5,7,8	0.551 –	0.671 7
<i>glass-0-1-2-3-vs-4-5-6</i>	0.988 7,8	0.989 7,8	0.894 8	0.989 7,8	0.879 8	0.984 7,8	0.731 –	0.722 –
<i>glass0</i>	0.937 3,5,7,8	0.929 3,5,7,8	0.792 5,7	0.930 3,5,7,8	0.724 –	0.881 5,7,8	0.628 –	0.675 –
<i>glass1</i>	0.913 3,5,6,7,8	0.918 3,5,6,7,8	0.681 –	0.830 5,7	0.668 –	0.790 7	0.546 –	0.669 –
<i>heart</i>	0.846 3,5,7,8	0.855 3,5,7,8	0.728 8	0.877 3,5,7,8	0.662 8	0.860 3,5,7,8	0.632 8	0.374 –
<i>jm1</i>	0.866 3,5,7,8	0.862 3,5,7,8	0.398 –	0.906 1,2,3,5,6,7,8	0.788 3,7,8	0.858 3,5,7,8	0.720 3,8	0.467 –
<i>kc1</i>	0.865 3,5,7,8	0.860 3,5,7,8	0.405 –	0.902 1,2,3,5,6,7,8	0.786 3,7,8	0.858 3,5,7,8	0.734 3,8	0.399 –
<i>kc2</i>	0.877 3,8	0.889 3,7,8	0.487 –	0.908 3,7,8	0.832 3	0.882 3,7,8	0.769 3	0.757 3
<i>kr-vs-k-three-vs-eleven</i>	1.000 3,5,7,8	1.000 3,5,7,8	0.981 –	1.000 3,5,7,8	0.987 7,8	1.000 3,5,7,8	0.978 –	0.968 –
<i>kr-vs-k-zero-one-vs-draw</i>	1.000 3,5,7,8	1.000 3,5,6,7,8	0.972 8	1.000 3,5,6,7,8	0.982 8	0.999 3,5,7,8	0.981 8	0.691 –
<i>page-blocks0</i>	0.994 3,5,7,8	0.994 3,5,7,8	0.850 –	0.999 1,2,3,5,6,7,8	0.962 3,7,8	0.992 3,5,7,8	0.913 –	0.734 –
<i>pima</i>	0.870 3,5,7,8	0.871 3,5,7,8	0.736 7,8	0.890 3,5,7,8	0.725 7,8	0.850 3,5,7,8	0.645 8	0.508 –
<i>segment0</i>	0.996 3,5,6,7,8	0.997 3,4,5,6,7,8	0.859 7,8	0.993 3,5,6,7,8	0.831 7,8	0.975 3,5,7,8	0.794 8	0.628 –
<i>shuttle-c0-vs-c4</i>	1.000 5,8	1.000 5,8	0.999 5,8	1.000 5,8	0.500 –	1.000 5,8	0.999 5,8	0.480 –
<i>vehicle0</i>	0.988 3,5,6,7,8	0.990 3,5,6,7,8	0.881 8	0.987 3,5,7,8	0.834 –	0.970 3,5,7,8	0.841 –	0.767 –
<i>vowel0</i>	0.997 3,5,6,7,8	0.999 3,5,6,7,8	0.852 8	0.999 3,5,6,7,8	0.862 8	0.989 3,5,7,8	0.870 8	0.260 –
<i>wisconsin</i>	0.998 7,8	0.997 7,8	0.929 8	0.997 7,8	0.716 8	0.995 7,8	0.958 8	0.081 –
<i>yeast-0-2-5-6-vs-3-7-8-9</i>	0.951 3,5,7,8	0.943 3,5,7,8	0.789 8	0.989 1,2,3,5,6,7,8	0.807 8	0.929 3,5,7,8	0.803 8	0.444 –
<i>yeast-0-2-5-7-9-vs-3-6-8</i>	0.983 3,5,7,8	0.979 3,5,7,8	0.896 8	0.997 2,3,5,6,7,8	0.925 8	0.978 3,5,7,8	0.900 8	0.252 –
<i>yeast-0-3-5-9-vs-7-8</i>	0.947 3,5,7,8	0.940 3,5,7,8	0.659 8	0.989 3,5,6,7,8	0.741 8	0.889 3,5,7,8	0.730 8	0.270 –
<i>yeast-0-5-6-7-9-vs-4</i>	0.984 3,5,7,8	0.975 3,5,7,8	0.867 8	0.994 3,5,6,7,8	0.891 8	0.965 3,5,7,8	0.828 8	0.089 –
<i>yeast-2-vs-4</i>	0.987 3,7,8	0.986 3,7,8	0.906 8	1.000 1,2,3,5,6,7,8	0.911 8	0.982 3,7,8	0.891 8	0.326 –
<i>yeast1</i>	0.887 3,5,7,8	0.885 3,5,7,8	0.768 7,8	0.922 3,5,7,8	0.754 7,8	0.877 3,5,7,8	0.655 8	0.339 –
<i>yeast3</i>	0.987 3,5,7,8	0.984 3,5,7,8	0.919 8	0.996 1,2,3,5,6,7,8	0.953 7,8	0.982 3,5,7,8	0.904 8	0.380 –

5 Conclusions and Future Directions

The proposed in this paper methods for imbalanced stream classification DSC-R and DSC-S performed favorably in comparison with other dedicated algorithms. The evaluation of the predictive abilities of the techniques was conducted on the basis of computer experiments. The obtained results were analyzed using statistical tests and for all the chosen metrics F-score, Gmean and AUC score, the proposed methods obtained satisfactory results, comparable to algorithms such as REA or Learn++.CDS or Learn++.NIE. The algorithm utilizes memory buffer in order to propagate the instances from the previous block that were chosen as the representatives. Since the buffer is of fixed size, after it is full some instances must be removed from it. In the current implementation, the oldest examples are deleted. A more advanced “forgetting” mechanism, that could favor the instances from the minority class and only the instances from the majority that are the best representatives could be introduced in order to further improve the performance of the classifier. Additionally testing other sampling methods for under- and over-sampling may prove to produce better results.

Acknowledgement. This work was supported by the Polish National Science Centre under the grant No. 2017/27/B/ST6/01325 as well as by the statutory funds of the Department of Systems and Computer Networks, Faculty of Electronics, Wrocław University of Science and Technology.

References

1. Alcalá-Fdez, J., Fernández, A., Luengo, J., Derrac, J., García, S.: Keel data-mining software tool: data set repository, integration of algorithms and experimental analysis framework. *Mult. Valued Log. Soft Comput.* **17**(2–3), 255–287 (2011). <http://dblp.uni-trier.de/db/journals/mvl/mvl17.html>
2. Bifet, A., Gavaldà, R.: Learning from time-changing data with adaptive windowing. In: *SIAM International Conference on Data Mining* (2007)
3. Bobowska, B., Woźniak, M.: Experimental study on modified radial-based over-sampling. In: Graña, M., et al. (eds.) *SOCO'18-CISIS'18-ICEUTE'18 2018. AISC*, vol. 771, pp. 110–119. Springer, Cham (2019). https://doi.org/10.1007/978-3-319-94120-2_11
4. Brzezinski, D., Stefanowski, J.: Ensemble classifiers for imbalanced and evolving data streams, pp. 44–68, March 2018. https://doi.org/10.1142/9789813228047_0003
5. Chen, S., He, H.: SERA: selectively recursive approach towards nonstationary imbalanced stream data mining. In: *2009 International Joint Conference on Neural Networks*, pp. 522–529, June 2009. <https://doi.org/10.1109/IJCNN.2009.5178874>
6. Chen, S., He, H.: Towards incremental learning of nonstationary imbalanced data stream: a multiple selectively recursive approach. *Evol. Syst.* **2**(1), 35–50 (2011). <https://doi.org/10.1007/s12530-010-9021-y>
7. Ditzler, G., Polikar, R.: Incremental learning of concept drift from streaming imbalanced data. *IEEE Trans. Knowl. Data Eng.* **25**(10), 2283–2301 (2013)

8. Ditzler, G., Roveri, M., Alippi, C., Polikar, R.: Learning in nonstationary environments: a survey. *IEEE Comput. Intell. Mag.* **10**, 12–25 (2015). <https://doi.org/10.1109/MCI.2015.2471196>
9. Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., Bouchachia, A.: A survey on concept drift adaptation. *ACM Comput. Surv.* **46**(4), 44:1–44:37 (2014). <https://doi.org/10.1145/2523813>. <http://doi.acm.org/10.1145/2523813>
10. Gama, J.: *Knowledge Discovery from Data Streams*, 1st edn. Chapman & Hall/CRC, Boca Raton (2010)
11. Gao, J., Ding, B., Fan, W., Han, J., Philip, S.Y.: Classifying data streams with skewed class distributions and concept drifts. *IEEE Internet Comput.* **12**(6), 37–49 (2008)
12. Hoens, T.R., Polikar, R., Chawla, N.V.: Learning from streaming data with concept drift and imbalance: an overview. *Prog. Artif. Intell.* **1**(1), 89–101 (2012). <https://doi.org/10.1007/s13748-011-0008-0>
13. Hulten, G., Spencer, L., Domingos, P.: Mining time-changing data streams. In: *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2001*, pp. 97–106. ACM, New York (2001). <https://doi.org/10.1145/502512.502529>. <http://doi.acm.org/10.1145/502512.502529>
14. Kolter, J.Z., Maloof, M.A.: Dynamic weighted majority: a new ensemble method for tracking concept drift. In: *Proceedings of the Third IEEE International Conference on Data Mining, ICDM 2003*, p. 123. IEEE Computer Society, Washington, D.C. (2003). <http://dl.acm.org/citation.cfm?id=951949.952136>
15. Koziarski, M., Krawczyk, B., Woźniak, M.: Radial-based approach to imbalanced data oversampling. In: Martínez de Pisón, F.J., Urraca, R., Quintián, H., Corchado, E. (eds.) *HAIS 2017. LNCS (LNAI)*, vol. 10334, pp. 318–327. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-59650-1_27
16. Krawczyk, B., Minku, L.L., Gama, J., Stefanowski, J., Woniak, M.: Ensemble learning for data stream analysis. *Inf. Fusion* **37**(C), 132–156 (2017). <https://doi.org/10.1016/j.inffus.2017.02.004>
17. Lemaître, G., Nogueira, F., Aridas, C.K.: Imbalanced-learn: a Python toolbox to tackle the curse of imbalanced datasets in machine learning. *J. Mach. Learn. Res.* **18**(17), 1–5 (2017). <http://jmlr.org/papers/v18/16-365.html>
18. Littlestone, N., Warmuth, M.K.: The weighted majority algorithm. *Inf. Comput.* **108**(2), 212–261 (1994). <https://doi.org/10.1006/inco.1994.1009>
19. Masud, M., Gao, J., Khan, L., Han, J., Thuraisingham, B.: A practical approach to classify evolving data streams: training with limited amount of labeled data, pp. 929–934, December 2008. <https://doi.org/10.1109/ICDM.2008.152>
20. Pedregosa, F., et al.: Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011). <http://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf>
21. Sayyad Shirabad, J., Menzies, T.: *The PROMISE repository of software engineering databases*. School of Information Technology and Engineering, University of Ottawa, Canada (2005). <http://promise.site.uottawa.ca/SERepository>
22. Stefanowski, J., Brzezinski, D.: Stream classification. In: Sammut, C., Webb, G.I. (eds.) *Encyclopedia of Machine Learning and Data Mining*, pp. 1191–1199. Springer, Boston (2017). https://doi.org/10.1007/978-1-4899-7687-1_908

23. Street, W.N., Kim, Y.: A streaming ensemble algorithm (SEA) for large-scale classification. In: Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2001, pp. 377–382. ACM, New York (2001). <https://doi.org/10.1145/502512.502568>. <http://doi.acm.org/10.1145/502512.502568>
24. Sun, Y., Wong, A.K.C., Kamel, M.S.: Classification of imbalanced data: a review. *IJPRAI* **23**, 687–719 (2009)
25. Wang, Y., Zhang, Y., Wang, Y.: Mining data streams with skewed distribution by static classifier ensemble. In: Chien, B.C., Hong, T.P. (eds.) Opportunities and Challenges for Next-Generation Applied Intelligence. SCI, vol. 214, pp. 65–71. Springer, Heidelberg (2009). https://doi.org/10.1007/978-3-540-92814-0_11
26. Woźniak, M., Graña, M., Corchado, E.: A survey of multiple classifier systems as hybrid systems. *Inf. Fusion* **16**, 3–17 (2014)