



**IJIRCCCE**

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

**Volume 10, Issue 6, June 2022**

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

**Impact Factor: 8.165**



9940 572 462



6381 907 438



ijircce@gmail.com



www.ijircce.com

# Radial-Based Oversampling for Multiclass Imbalanced Data Classification

Dr. C. M. Suvarna varma<sup>1</sup>, Muddana Sravani<sup>2</sup>, Koutharapu Mohana<sup>3</sup>, Movva Mounvitha<sup>4</sup>

Assoc. Professor, Department of IT, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur,  
Andhra Pradesh, India<sup>1</sup>

UG Student, Department of IT, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur, Andhra Pradesh, India<sup>2</sup>

UG Student, Department of IT, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur, Andhra Pradesh, India<sup>3</sup>

UG Student, Department of IT, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur, Andhra Pradesh, India<sup>4</sup>

**ABSTRACT:** Learning from imbalanced data is among the most popular topics in the contemporary machine learning. However, the vast majority of attention in this field is given to binary problems, while their much more difficult multiclass counterparts are relatively unexplored. Handling data sets with multiple skewed classes poses various challenges and calls for a better understanding of the relationship among classes. It proposes multiclass radial-based oversampling (MC-RBO), a novel data-sampling algorithm dedicated to multiclass problems. The main novelty of the method lies in using potential functions for generating artificial instances. It takes into account information coming from all of the classes, contrary to existing multiclass oversampling approaches that use only minority class characteristics.

The process of artificial instance generation is guided by exploring areas where the value of the mutual class distribution is very small. This way, it ensure a smart oversampling procedure that can cope with difficult data distributions and alleviate the shortcomings of existing methods. The usefulness of the MC-RBO algorithm is evaluated on the basis of extensive experimental study and backed-up with a thorough statistical analysis. Obtained results show that by taking into account information coming from all of the classes and conducting a smart oversampling, It can significantly improve the process of learning from multiclass imbalanced data.

## I. INTRODUCTION

Due to its compatibility with real-world pattern classification problems, where the most significant or interesting classes are typically strongly underrepresented, learning from imbalanced data has become the subject of intense research. The bulk of research concentrate on the binary problem, in which the well-represented class is referred to as the majority, while the underrepresented class is referred to as the minority. Learning difficulties is common in this situation, as standard classification algorithms are biased toward the majority class. However, in the situation of unbalanced data classification, learning challenges are caused by more than just the disparity across classes. Despite a high imbalance ratio, one can easily come up with an example where the instance distributions from distinct classes are well separated (IR).

Solutions at the Data Level: These methods rebalance the training set such that it may be used by any conventional classifier. This is accomplished by either reducing the number of instances of the majority class (under sampling) or increasing the number of examples of the minority class (oversampling) (oversampling). Because random procedures provide no control over the rebalancing process, guided solutions that retain minority class features have been investigated. Synthetic Minority Oversampling Technique (SMOTE), which is now considered a cornerstone for the bulk of suggested oversampling algorithms, is the most prominent approach based on this philosophy. Unfortunately, because SMOTE implies that instances from the minority class form groups, it may result in changes in the minority class's characteristics and, as a result, over fitting the classifier.

## MOTIVATION:

Algorithm-level approaches aim at modifying the classifier learning procedure in order to make it skew-insensitive.

This requires an in-depth understanding of the modified methods, as well as of the actual learning difficulty that causes the poor performance on minority class. Here, cost-sensitive approaches are popular, as they allow to easily modify any learning method by adding a separate misclassification penalty for each class [8]. This should improve minority class recognition, as classifier will be much more penalized for misclassification of minority instance. Another potential solution include usage of one-class classifiers [3]. Here, we create a data description of the target class (one selected by the user) and treat the remaining one as outliers. While we sacrifice knowledge about one of the classes, we gain a skew-insensitive classifier that captures unique properties of its target. Hybrid solutions use advantages of the mentioned approaches and combine them with other methodologies, mainly ensemble learners [16]. They take advantage of increased predictive power, diversity and ability to capture complex data offered by combined classifiers and augment it with tackling imbalance at the level of each classifier. Popular approaches include combination of Bagging or Boosting with pre-processing.

#### AIM AND SCOPE:

We present a novel oversampling technique that employs radial functions to estimate the potential of occurrences in this work. We propose that they be used to simulate mutual class distributions and assess the learning difficulty associated with each instance. Instead of using all of the objects, our technique allows you to choose which ones should be oversampled. We can predefine the nature of minority class instances and use it to direct the artificial instance injection procedure by examining differences in potential at a specific point. This allows for a more accurate representation of the underlying distribution of the minority class. Our method is also useful for applications in high-dimensional datasets because it does not rely on neighbourhood calculations. Experiments on a variety of benchmarks show that the suggested radial-based oversampling is capable of delivering satisfactory results.

#### GOALS AND OBJECTIVES:

- Proposition of the MC-RBO algorithm, which allows for intelligent data oversampling that exploits local data characteristics of each class.
- New approach for generating artificial instances that do not use NNs of each minority instance for instance imputation, thus making MC-RBO robust to a typical data distributions.

## II. LITERATURE REVIEW

**B. Krawczyk [1]**, Multi-class imbalanced classification is not as well-developed as its binary counterpart. Here we deal with a more complicated situation, as the relations among the classes are no longer obvious. A class may be a majority one when it is compared to some other classes, but a minority or well-balanced for the rest of them. When dealing with multi-class imbalanced data we may easily lose performance on one class while trying to gain it on another. Considering this problem there are many issues that must be addressed by novel proposals. A deeper insight into the nature of the class imbalance problem is needed, as one should know in what domains does class imbalance most hinder the performance of standard multi-class classifiers when designing a method tailored for this problem. While most of the challenges discussed can also be transferred to multi-class problems, there is a number of topics specific just to them. We identify the following vital future directions of research.

**P. Branco, L. Torgo, and R. P. Ribeiro [2]**, Many real-world data-mining applications involve obtaining predictive models using datasets with strongly imbalanced distributions of the target variable. Frequently, the least-common values of this target variable are associated with events that are highly relevant for end users (e.g., fraud detection, unusual returns on stock markets, anticipation of catastrophes, etc.). Moreover, the events may have different costs and benefits, which, when associated with the rarity of some of them on the available training data, creates serious problems to predictive modeling techniques. This article presents a survey of existing techniques for handling these important applications of predictive analytics. Although most of the existing work addresses classification tasks (nominal target variables), we also describe methods designed to handle similar



problems within regression tasks (numeric target variables). In this survey, we discuss the main challenges raised by imbalanced domains, propose a definition of the problem, describe the main approaches to these tasks, propose a taxonomy of the methods, summarize the conclusions of existing comparative studies as well as some theoretical analyses of some methods, and refer to some related problems within predictive modeling.

**K. Napierala and J. Stefanowski [3]**, A dataset is considered to be imbalanced if it is characterized by an unequal distribution of examples between the classes. The smaller class is called the minority class, while the other classes are called majority classes. In the imbalanced datasets, the minority class is usually of primary interest to the decision maker, i.e. not recognizing the minority class examples is much more serious than raising so called false alarms (assigning a majority example to the minority class). For this reason, the most popular performance measures such as total accuracy are not useful in the context of class imbalance, as they are biased towards the majority classes. Thus, for imbalanced domains other performance measures have to be used. The imbalance of a learning dataset can be either intrinsic (in the sense that it is a direct result of the nature of the data space) or it can be caused by too high costs of progress.

**N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer [5]**, An approach to the construction of classifiers from imbalanced datasets is described. A dataset is imbalanced if the classification categories are not approximately equally represented. Often real-world data sets are predominately composed of "normal" examples with only a small percentage of "abnormal" or "interesting" examples. It is also the case that the cost of misclassifying an abnormal (interesting) example as a normal example is often much higher than the cost of the reverse error. Under-sampling of the majority (normal) class has been proposed as a good means of increasing the sensitivity of a classifier to the minority class. This paper shows that a combination of our method of over-sampling the minority (abnormal) class and under-sampling the majority (normal) class can achieve better classifier performance (in ROC space) than only under-sampling the majority class. This paper also shows that a combination of our method of over-sampling the minority class and under-sampling the majority class can achieve better classifier performance (in ROC space) than varying the loss ratios in Ripper or class priors in Naive Bayes. Our method of over-sampling the minority class involves creating synthetic minority class examples. Experiments are performed using C4.5, Ripper and a Naive Bayes classifier. The method is evaluated using the area under the Receiver Operating Characteristic curve (AUC) and the ROC convex hull strategy.

**C. Bunkhumpornpat, K. Sinapiromsaran, and C. Lursinsap [6]**, Multiclass classification in cancer diagnostics using DNA or Gene Expression Signatures, but also classification of bacteria species fingerprints in MALDI-TOF mass spectrometry data, is challenging because of imbalanced data and the high number of dimensions with respect to the number of instances. In this study, a new oversampling technique called LICIC will be presented as a valuable instrument in countering both class imbalance, and the famous-curse of dimensionality problem. The method enables preservation of non-linearities within the dataset, while creating new instances without adding noise. The method will be compared with other oversampling methods, such as Random Oversampling, SMOTE, Borderline-SMOTE, and ADASYN. F1 scores show the validity of this new technique when used with imbalanced, multiclass, and high-dimensional datasets.

### III. SYSTEM ANALYSIS

#### 3.1 REQUIREMENT ANALYSIS:

Requirements analysis, also called requirements engineering, is the process of determining user expectations for a new or modified product. These features, called requirements, must be quantifiable, relevant and detailed. In software engineering, such requirements are often called functional specifications.

Requirements analysis is critical to the success or failure of a systems or software project. The requirements should be

documented, actionable, measurable, testable, traceable, related to identified business needs or opportunities, and defined to a level of detail sufficient for system design.

**User requirements:**

1. Execution time should be fast
2. More Accurate.

**Software requirements:**

1. Operating System: Windows 7/8/10.
2. Language: Python.
3. Data Base: FER2013.
4. Tools: Google Collabatory / Jupyter notebook

**Hardware requirements:**

1. SYSTEM: Intel Core i5-2600.
2. HARD DISK: 40 GB.
3. MONITOR: 15 VGA COLOUR.
4. MOUSE: Logitech.
5. RAM: 4 GB or more.

**Functional requirements:**

**Input:**

Defined MC-RBO and B-RBO classes which are used to oversample the datasets. Tests the data sets over different built in oversampling models.

And tests the accuracy scores of the data set using Random Forest Classifier. Compares the accuracy scores of different oversampled methods.

**Output:**

Output consists of showing that MCRBO is used to oversample better over the remaining methods for multi class imbalanced data classification and produce a great accuracy score when compared to other methods when tested using a ensemble classifier namely Random Forest Classifier. by comparing the methods accuracy score the output is shown



**Non Functional requirements:**

**Execution qualities:**

**Efficiency:**

The state or quality of being efficient, i.e., able to accomplish something with the least waste of time and effort; competency in performance.

**Evolution qualities:**

**Testability:**

The means by which the presence, quality, or genuineness of anything is determined

**Extensibility:**

To enlarge the scope of, or make more comprehensive, as operations, influence etc.

**Scalability:**

The ability of something, especially a computer system, to adapt to increased demands.

**Usability:**

The system is designed with completely automated process hence there is no or less user intervention.

**Reliability:**

The system is more reliable because of the qualities that are inherited from the chosen platform. The code built by using Python Programming is more reliable.

**Performance:**

This system is developing in the high-level languages and using the advanced front- end and back-end technologies it will give response to the end user on client system with in very less time.

**Supportability:**

The system is designed to be the cross platform supportable. The system is supported on a wide range of hardware and any software platform, which is having Python, built into the system.

**Implementation:**

The system is implemented in Python environment.

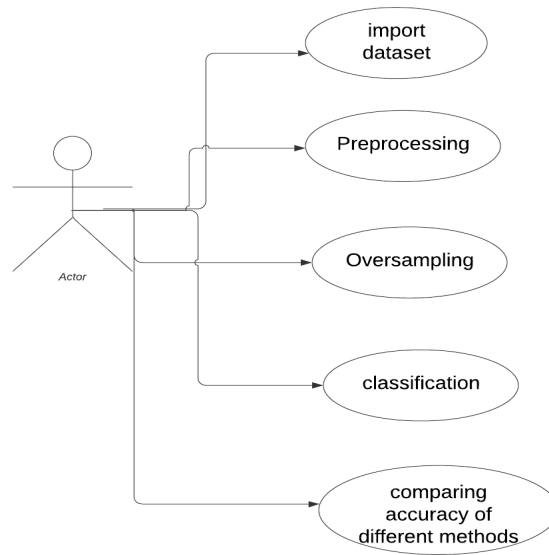


Fig: Use case Diagram

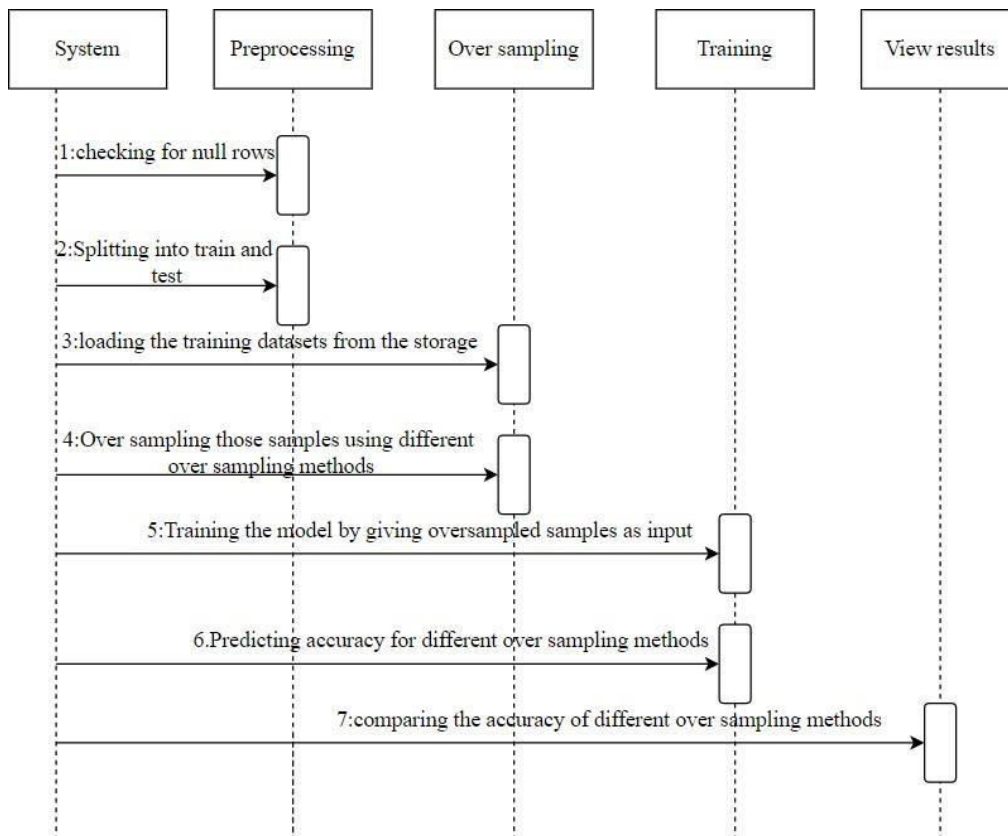


Fig: Sequence Diagram

#### **IV. EXISTING SYSTEM**

##### **Random OverSampling:**

Random resampling provides a naive technique for rebalancing the class distribution for an imbalanced dataset. Random oversampling duplicates examples from the minority class in the training dataset and can result in over fitting for some models. Random under sampling deletes examples from the majority class and can result in losing information invaluable to a model. Random oversampling involves randomly duplicating examples from the minority class and adding them to the training dataset. Examples from the training dataset are selected randomly with replacement. This means that examples from the minority class can be chosen and added to the new more balanced training dataset multiple times; they are selected from the original training dataset, added to the new training dataset, and then returned or -replaced| in the original dataset, allowing them to be selected again.

##### **Smooth Oversampling:**

SMOTE works by selecting examples that are close in the feature space, drawing a line between the examples in the feature space and drawing a new sample at a point along that line. Specifically, a random example from the minority class is first chosen. Then  $k$  of the nearest neighbours for that example are found (typically  $k=5$ ). A randomly selected neighbour is chosen and a synthetic example is created at a randomly selected point between the two examples in feature space. This procedure can be used to create as many synthetic examples for the minority class as are required. As described in the paper, it suggests first using random undersampling to trim the number of examples in the majority class, then use SMOTE to oversample the minority class to balance the class distribution. The approach is effective because new synthetic examples from the minority class are created that are plausible, that is, are relatively close in feature space to existing examples from the minority class.

##### **ADASYN OverSampling:**

The essential idea of ADASYN is to use a weighted distribution for different minority class examples according to their level of difficulty in learning, where more synthetic data is generated for minority class examples that are harder to learn compared to those minority examples that are easier to learn.

As a result, the ADASYN approach improves learning with respect to the data distributions in two ways

- (1) Reducing the bias introduced by the class imbalance, and
- (2) Adaptively shifting the classification decision boundary toward the difficult examples.

#### **IV. PROPOSED SYSTEM**

Handling data sets with multiple skewed classes poses various challenges and calls for a better understanding of the relationship among classes. In this paper, we propose multiclass radial-based oversampling (MCRBO), a novel data-sampling algorithm dedicated to multiclass problems. The main novelty of our method lies in using potential functions for generating artificial instances. We take into account information coming from all of the classes, contrary to existing multiclass oversampling approaches that use only minority class characteristics. The process of artificial instance generation is guided by exploring areas where the value of the mutual class distribution is very small. This way, we ensure a smart oversampling procedure that can cope with difficult data distributions and alleviate the shortcomings of existing methods.

#### **V. RESULT & DESCRIPTION**

Several measures were explored in order to appropriately analyse the behaviour of the examined algorithms on imbalanced data: accuracy, precision, recall, F-measure, and geometric mean (G- mean). Friedman test was used to ensure statistical validity of the results, and average rankings across all datasets were presented. The performance of the oversampling strategy suggested in this research was comparable to that of the SMOTE method. In terms of recall, using ADASYN produced the greatest results on the datasets studied, at the cost of slightly reduced precision.



Dataset	MC-RBO	ADASYN	SMOTE	Random Over sampling
Automobile	0.72	0.75	0.76	0.72
Balance	0.79	0.80	0.80	0.82
Car	0.96	0.95	0.95	0.96
Cleveland	0.55	0.56	0.55	0.58
Contraceptive	0.54	0.54	0.54	0.54
Dermatology	0.98	0.94	0.98	0.98
Ecoli	0.83	0.82	0.82	0.77
Flare	0.71	0.73	0.74	0.74
Glass	0.68	0.71	0.70	0.71
Hayes-Roth	0.74	0.80	0.72	0.72
Led7digit	0.64	0.64	0.67	0.65
Lymphography	0.86	0.87	0.83	0.86
New-thyroid	0.91	0.89	0.91	0.91
Page-blocks	0.96	0.95	0.96	0.97
Thyroid	0.98	0.97	0.97	0.98
Vehicle	0.75	0.74	0.74	0.75
Wine	0.96	0.95	0.96	0.96
Wine quality red	0.63	0.64	0.61	0.62
Yeast	0.55	0.57	0.55	0.56
Zoo	0.96	0.96	0.96	0.94

## VI. CONCLUSION

The main aim was to propose a novel, effective oversampling approach for a challenging task of multiclass imbalanced data classification, proposed the MC-RBO algorithm, which can exploit local data characteristics of each class and offers better placement of new artificially generated instances and more targeted empowering of minority classes. As distinct from SMOTE-like algorithms, MC-

RBO does not use NNs of minority instances, and thus, it is robust to atypical data distributions, especially when the minority classes form several disjoint clusters. The computer experiments confirmed the usefulness of the proposed approach and on the basis of a thorough statistical analysis we may assert that for many types of data sets MC-RBO is statistically significantly better than state-of-art methods. Particularly, it outperforms existing multiclass oversampling methods, as well as typical binary decomposition scenarios.

## REFERENCES

1. B. Krawczyk, —Learning from imbalanced data: Open challenges and future directions,| Prog. Artif. Intell., vol. 5, no. 4, pp. 221–232, Nov. 2016.
2. P. Branco, L. Torgo, and R. P. Ribeiro, —A survey of predictive modeling on imbalanced domains,| ACM Comput. Surv., vol. 49, no. 2, p. 31, Nov. 2016.
3. K. Napierala and J. Stefanowski, —Identification of different types of minority class examples in imbalanced data,| in Proc. Int. Conf. Hybrid Artif. Intell. Syst., Mar. 2012, pp. 139–150.
4. X.-W. Chen and M. Wasikowski, —FAST: A roc-based feature selection metric for small samples and imbalanced



data classification problems,|| in Proc. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, Aug. 2008, pp. 124–132.

5. N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, —SMOTE: Synthetic minority over-sampling technique,|| J. Artif. Intell. Res., vol. 16, no. 1, pp. 321–357, 2002.
6. C. Bunkhumpornpat, K. Sinapiromsaran, and C. Lursinsap, —Safelevel-smote: Safe-level- synthetic minority over-sampling technique for handling the class imbalanced problem,|| in Proc. Asia–Pacific Conf. Knowl. Discovery Data Mining, Apr. 2009, pp. 475–482.
7. M. Kubat and S. Matwin, —Addressing the curse of imbalanced training sets: One-sided selection,|| in Proc. 14th Int. Conf. Mach. Learn., Jul. 1997, pp. 179–186.
8. S. Wang and X. Yao, —Multiclass imbalance problems: Analysis and potential solutions,|| IEEE Trans. Syst., Man, Cybern. B, Cybern., vol. 42, no. 4, pp. 1119–1130, Aug. 2012.
9. A. Fernández, S. García, M. Galar, R. C. Prati, B. Krawczyk, and F. Herrera, Learning From Imbalanced Data Sets. Cham, Switzerland: Springer, 2018.
10. J. A. Sáez, B. Krawczyk, and M. Wozniak, —Analyzing the oversampling of different classes and types of examples in multi-class imbalanced datasets,|| Pattern Recognit., vol. 57, pp. 164–178, Sep. 2016.



INNO  SPACE  
SJIF Scientific Journal Impact Factor

Impact Factor: 8.165

 **doi**<sup>®</sup>  
**cross** **ref**

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 9940 572 462  6381 907 438  [ijircce@gmail.com](mailto:ijircce@gmail.com)



[www.ijircce.com](http://www.ijircce.com)

Scan to save the contact details