

*Original Article*

## Robust Machine Learning Models for Imbalanced Dataset Classification

**Dr. Vasanth Kumar<sup>1</sup>, Girija Rajan<sup>2</sup>**

<sup>1,2</sup>Department of Artificial Intelligence and Data Science, Arjun College of Technology, Pollachi, Tamil Nadu, India.

Received: 03-12-2025

Revised: 25-12-2025

Accepted: 01-01-2026

Published: 07-01-2026

### ABSTRACT

The problem of class imbalance in machine learning classification is widely present and difficult across the machine learning area, especially in real-world tasks, including fraud detection, medical diagnosing, network intrusion detection and fault prediction. When this occurs, the minority population is more likely to capture the important occurrences and the conventional machine learning models normally focus on the majority population and give misleading accuracy with poor generalization and high costs of misclassification. This paper is the result of an extensive research into powerful machine learning techniques in the classification of imbalanced datasets. The paper presents a systematic review of theoretical underpinnings of learning imbalance, literature reviews on state-of-the-art methods, such as data, algorithm-level and ensemble based methods, and suggests a convergent system methodology to build a robust classifier. Linear resampling algorithms, cost-effective learning algorithms, hybrid ensemble algorithms, and imbalanced data evaluation metrics are discussed in details. An organized experimental procedure is described to measure robustness when imbalance ratios and various noise levels are changing. Comparative findings indicate that hybrid methods that combine adaptive resampling and cost sensitive loss functions are always better than simpler classifiers based on their F1-score, G-mean, and area under the precision-recall curve. The discussion demonstrates practical trade-offs between model performance, model complexity and interpretability. In the conclusion part, the paper highlights future research directions which include scalable imbalance learning, deep learning adaptations and domain aware evaluation strategies. The paper is an excellent source of information to a researcher and practitioner aiming at finding principled and effective solutions to imbalanced classification problems.

### KEYWORDS

Imbalanced Datasets, Robust Classification, Cost-Sensitive Learning, Ensemble Learning, Resampling Techniques, Minority Class Prediction, Machine Learning Evaluation Metrics.

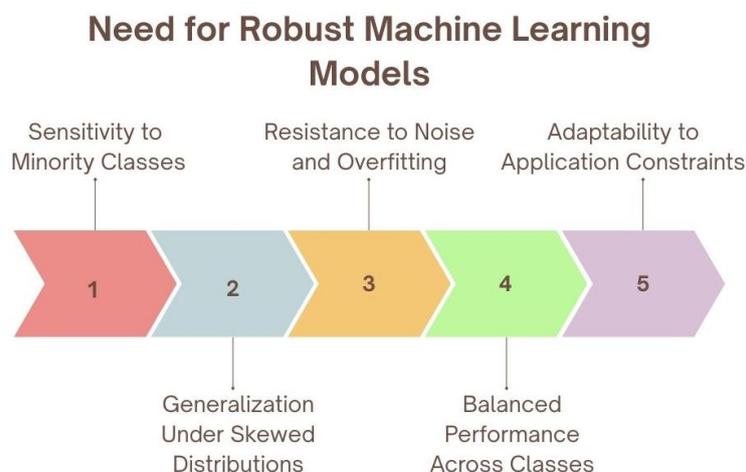
## 1. INTRODUCTION

### 1.1. Background

Imbalanced information situations are intrinsic to most of the areas of application of high significance in the real world where the most important events are infrequent but have substantial impacts. Fraudulent transactions in financial systems are just a minor percentage among the entire transactions that have been registered, but the consequence of inability to detect them may cause massive monetary losses. The same thing can be said about medical diagnostics where malignant tumors are much rarer in comparison with benign ones, however, early and precise diagnosis is essential to save the life of a patient. Failure and abnormalities in industrial and cyber-physical systems are infrequent occurrences, but these situations must be properly detected before disastrous consequences appear. These examples point to one of the main problems, which is that, though there are not many cases of ethnic minorities classes, they usually count a lot more than those of the majority class, and their effective identification should be a priority. According to the statistical learning theory, class imbalance creates serious bias in the empirical fitting of the risk minimization process that forms the backbone of the majority of supervised learning algorithms. The errors are uniformly distributed across all the samples in standard loss functions, so such loss functions result in majority-class samples overwhelming the optimisation goal. This leads to skewed learned boundaries to majority parts of the feature space, and leaves minority parts of the feature space poorly modelled or even ignored. With this, the traditional classifiers like logistic regression, decision trees and support vector machine would typically attain good overall accuracy and very low sensitivity or recall of the minority class. This effect makes such models virtually useless within imbalance-sensitive processes, where the accurate classification of rare events is much more useful than accurate classification of normal cases of high abundance. These constraints are a very strong incentive to designing special learning instructions that clearly consider the imbalance in the classes and focuses on minority-class detection at the cost of generalization performance.

### 1.2. Need for Robust Machine Learning Models

Effective task of imbalanced dataset handling requires to have strong machine learning models since the traditional learning algorithms do not perform well in case of a skewed distribution of classes. Several practical as well as theoretical reasons necessitate the requirement of robustness which are as follows.



**Fig 1 - Need for Robust Machine Learning Models**

### 1.2.1. Sensitivity to Minority Classes

In the case of imbalanced learning, the minority group of learning tends to represent very unusual but important events. Strong models should be sensitive to such minority cases so that they are picked correctly. This needs learning patterns that do not accommodate the overbearing of the minority group by the majority group in training.

### 1.2.2. Generalization under Skewed Distributions

Regular model must have the same or different ratios of imbalances of unknown data to remain stable at exposure of unknown data. The fact that some minorities are overfitted on oversampled minorities than others or a high rate of underfitting by the majority can be disastrous. Consequently, robustness would mean the capacity to make generalization between the imbalance levels without a severe decline in predictive accuracy.

### 1.2.3. Resistance to Noise and Overfitting

The samples that belong to minority-class are periodic and can be noisy. Strong models have to take these limited examples into account without being overwhelmed by them, and also be able to produce meaningful patterns. This is a balance that is crucial to have bound reliable decision boundaries and more so in high dimensional feature spaces.

### 1.2.4. Balanced Performance Across Classes

Instead of unbiased accuracy maximization, the goal of robust models is to have a balanced performance by all classes. This incorporates reasonable precision, recall and stability of both minority and majority classes which is not only fair but effective in real-life applications.

### 1.2.5. Adaptability to Application Constraints

Various fields have varying requirements including interpretability, computational efficiency, or real-time prediction. Strong machine learning models ought to be adaptable enough to align with such restrictions and at the same time perform reasonably well in class imbalance.

## 1.3. Models for Imbalanced Dataset Classification

Imbalanced dataset classification models are targeted to address the drawback inherent in a traditional learning algorithm that focuses on the largest classes of data elements to the disadvantage of all others: they present imbalance-knowledge during both training and decision-making. The conventional classifiers that are logistic regression, decision trees, support vector machines and neural networks usually optimize aggregate accuracy or likelihood-based goals which are mainly dominated by majority-class samples when class distributions are skewed. Consequently, such models tend to be poor recallers of minority classes although they may have high overall accuracy. In order to overcome this weakness, various specialization modeling strategies have been designed that vary the representation of data, the learning goal or the structure of the model itself. The first method that has been popularly used is the idea of improving traditional models with cost-sensitive learning where class-specific meat loss is incorporated in the loss function. Working with models that create decision boundaries that are more sensitive to minority-class example errors can be motivated by increasing the penalty on such errors (such as cost-sensitive logistic regression, support vector machines, and neural networks). A different type of models, which uses ensemble learning, encompass balanced random forests, boosting-based models, and hybrid ensembles, based on a combination of multiple base learners, where each is trained on a balanced or resampled dataset. The inclusion of so much diversity and robustness benefits of these ensemble models result in better

minority-class detection and variations. In more recent work, deep learning models as well are extended to imbalanced classification using special loss functions, including weighted cross-entropy and focal loss, which focus on the hard-to-classify samples and minority samples. Also, hybrid frameworks combining resampling, cost-sensitive optimization, and ensemble aggregation have been shown to be very effective in a large variety of imbalance conditions. Taken together, these models indicate the trend of the replacement of the accuracy-based learning by the balanced and application-sensitive classification, which is why the given models could be discussed as a great necessity in the context of the efficient work with imbalanced datasets in the framework of machine learning applications.

## 2. LITERATURE SURVEY

### 2.1. Data-Level Approaches

Data level techniques attempt to reduce the class imbalance by directly manipulating the training data prior to the application of a learning algorithm. These methods can be described as model-agnostic, that is, they do not change anything inside the classifier. Data-level methods can, by either sampling, or generating data, balance the distributions of classes in an effort to supply learning algorithms with a better representation of the minority class, which enhances classification performance.

#### 2.1.1. *Random Oversampling and Undersampling*

One of the most basic and oldest methods of managing class imbalance is random oversampling and undersampling. Random oversampling randomly repeats the existing samples of minority-class samples to augment the samples, which may assist the classifiers in learning more about the decision regions of minorities. Nonetheless, this overfitting can be frequent in such a duplication because the model can remember repeat samples, instead of making the generalization. Conversely, random undersampling lowers the sample sizes of the majority-class, which can be much faster to train and less biased against the majority class. However, in informative samples, undersampling can exclude important samples, which can worsen the overall classification, particularly in the case of a small original population.

#### 2.1.2. *Synthetic Sample Generation*

Simple duplication has its limitations, which have been addressed with synthetic sample generation methods which generate new samples of the minority-class by interpolating existing minority samples in the feature space. These approaches will add more diversity to the minority-class without altering its structure. This is done by producing synthetic data instead of copying samples to minimize overfitting and to better learn more accurate and smooth decision boundaries by the classifier. It has been demonstrated that synthetic sampling is especially useful in high-dimensional spaces where the samples of the minority are sparse and ill represented.

### 2.2. Algorithm-Level Approaches

The approaches at the algorithm level correct the issue of class imbalance by making adjustments to the learning process instead of the input data. These approaches directly include imbalance awareness in the model training process, typically by modifying loss functions, decision policies or optimization policies. In contrast to data-level approaches, algorithm-level approaches are generally classifier-specific, but potentially give principled and more theoretically-founded solutions to imbalance issues.

### 2.2.1. Cost-Sensitive Learning

Cost-sensitive learning explicitly takes into account various misclassification costs incurred on various classes with greater penalties on those incurred on minority-class instances. This strategy changes the learning goal such that the penalty associated with misclassifying majority sample is reduced. Consequently, this induces the classifier to pay more attention to correctly classifying the minority ones even at the cost of some degree in the majority-class performance. Cost sensitive schemes are most useful where the misclassification costs are known or estimable, e.g when medical or financial miscollected information is involved and false negatives are more expensive than false positives.

### 2.2.2. Threshold Adjustment

Threshold adjustment is a post-processing mechanism that changes the decision threshold applied to give out class labels and particularly probabilistic classifiers. A default threshold (0.5 in binary classification) is not used but instead the threshold is biased towards minority-class prediction. This method does not involve retraining of the model and it is computationally efficient. Practitioners can directly manipulate the trade-off between precision and recall by changing thresholds, thus the technique is particularly helpful in those applications where recall of the minority classification is a necessary requirement.

## 2.3. Ensemble-Based Approaches

Ensemble-based methods utilize several base learners in order to enhance predictive factors and resilience when classes are imbalanced. Ensemble methods reduce both variance and bias by averaging a variety of models or training on dissimilar balanced subsets of data. These methods have gained great popularity because of their good empirical effectiveness in a broad spectrum of imbalanced learning tasks.

### 2.3.1. Bagging-Based Ensembles

Ensemble methods using bagging solve the issue of class imbalance by creating, using a bootstrap, several balanced class samples. The learner statistics of each of the bases are trained on a balanced different subset, which enables the ensemble to learn the minority-class pattern better. The balanced bagging averages predictions made by the models which minimize the variance and enhances the recall of minority classes. This strategy has been found to work particularly well with unstable learners, including decision trees which are highly advantaged by bootstrap aggregation.

### 2.3.2. Boosting for Imbalanced Data

Boosting algorithms are used to improve the performance of the weak learners by focusing on those misclassified in the earlier iterations. With this scenario of imbalanced data, imbalance-sensitive boosting methods adjust the weight update algorithm to emphasize more on the sample of the minority classes. This will allow minority cases to get more and more attention during the training process. Consequently, one can train boosting-based methods to acquire decision boundaries of complex forms and has shown excellent performance in skewed datasets, especially when used with cost-sensitive or sampling based methods.

## 3. METHODOLOGY

### 3.1. Problem Formulation

In most practical classification problems, the data available has a highly skewed distribution of the class between the instances of one class being significantly greater than the other. Assume the

dataset is to be modeled as the set of input output pairs with an input being an input feature vector and an output being the binary label of a class. The label of the classes is that of one, or an integer of zero, where the minority class consists of a label of one and the majority category consists of a label of zero. One of the characteristic features of such datasets is that the likelihood of the occurrence of an instance of the minority-class is significantly smaller in comparison with the likelihood of the occurrence of an instance of the majority-class. This asymmetry confronts the traditional learning algorithms which are generally inclined to maximize total accuracy and thus incline towards a majority group with a result that the minority examples are poorly detected. The main goal within this problem context, is to acquire a classification system which by inputting feature vectors on features of input, classifies a given data point to its appropriate category within a system in a manner that effectively labels minority-class instances without only compromising on the generalization ability in prediction on unseen data. Whereas in a normal classification problem, it suffices to maximize overall accuracy, in a classification problem where produced accuracy is desirable, a classifier may optimize overall accuracy by simply relying on the majority of the investment and predicting its most common class. Rather, the learning process should focus on providing the appropriate identification of minority-class samples, which can usually be the most crucial in practice (fraud detection, disease diagnosis and fault detection). Meanwhile, the classifier should not be overfitted to the small minority samples such that the learned decision boundaries are strong and cross-generalizable. This involves establishing a delicate trade between sensitive to the minority class and stable in the whole space of features. The problem formulation is, therefore, through the design or selection of learning strategies that effectively deal with the problem of imbalance in classes directly, whether through dispersing the data, altering the learning goal, or a combination of the models, such that minority-class identification in the sample is improved to the general predictive reliability level.

### 3.2. Proposed Framework

The proposed structure is aimed at successfully dealing with the issues provided by the inequality in the number of classes integrating complementary approaches at various phases of learning. It combines adaptive resampling technique, cost-sensitive learning, and ensemble aggregation that enhances the capability to detect the minority-class and at the same time exhibits a strong generalization ability.

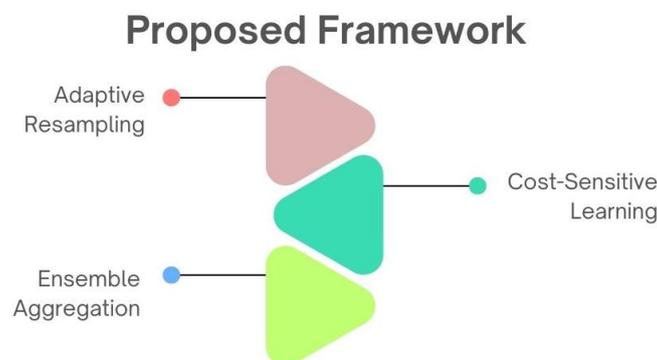


Fig 2 - Proposed Framework

#### 3.2.1. Adaptive Resampling

The idea behind adaptive resampling is to adaptively change the distribution of classes of the training data in terms of data properties instead of using uniform sampling. This method is used to

upsample and downsample minority-class samples selectively in challenging or sparsely sampled parts of the feature space, unlike either random undersampling or oversampling. Adaptive resampling optimizes the balance of classes and better represents the minority patterns by concentrating on the informative minority examples and eliminating overfitting and information loss by reducing redundancy.

### 3.2.2. Cost-Sensitive Learning

Cost-sensitive learning directly includes the misclassification costs that affect a particular class into the training of the classifier to capture the unequal impact of the classification mistakes. The suggested system attaches a greater punishment to the misclassification of the cases of minority-classes in order to motivate the classifier to focus on their accurate prediction. This change in the learning objective can be seen to better align with the imbalanced information situation in both minority and equal recall and performance, and the model becomes more appropriate in that case.

### 3.2.3. Ensemble Aggregation

Ensemble aggregation uses a number of individual base classifiers trained to the sampling distribution or cost configuration anticipated to be minimally predictive to produce a more robust predictor. The framework lessens the variation and bias against the majority class through the summation of the contributions of various learners. This combination approach enhances the discrimination of the minority-classes and offers more stable forecasts, relying on the advantages of each of the models and countering their drawbacks.

## 3.3. Adaptive Resampling Strategy

Adaptive resampling is used to overcome the shortcomings of the conventional fixed-ratio sampling methods by controlling the resampling procedure in response to the underlying data distribution. This is particularly true in highly gendered datasets where the instances of the minority-class can be very sparse, and can even intersect with the majority-class areas, rendering homogenous oversampling ineffective and potentially detrimental. Instead of using a fixed value of the ratio of the over or undersampling, the suggested adaptive resampling approach examines the separability of classes and the density of local data to inform the training sample creation and screening. This enables the resampling procedure to target parts on the feature space that contains small numbers of minority samples or are challenging to categorize. In particular, the areas with ambiguity in the decision-making step dividing majority and minority classes are revealed by class separability measures. Sample of minorities that are found close to such boundaries are accorded greater significance because they are critical in establishing correct classification boundaries. When such samples are adaptively expanded, the representation of such samples allows the classifier to acquire finer nuances between classes. Simultaneously, the local density information is used to prevent unnecessarily sampling dense minority areas that may result in redundancy and overfitting. Seldom populated areas by minorities, conversely, are concentrated in so that they enhance coverage and minimize bias. Moreover, the adaptive approach permits the sparse poor in areas with clear separation of the classes so as to reduce the dominance whilst sparing informative data not necessarily when class separation is evident. This pionic treatment of both classes will help to make the resampled data representative of the original data distribution and remove effects of imbalance. Generally, adaptive resampling offers an information-driven adaptable structure that improves minority-class representation, sustains important structural data, and enables learning of robust and generalizable classifiers in skewed classification.

### 3.4. Cost-Sensitive Loss Design

In order to directly tackle the bias caused by the imbalance in classes when training a model, a weighted cross-entropy cost-sensitive loss function is used. In normal cross-entropy loss, every instance of misclassification is equal and this leads to learning algorithms focusing on the majority class in case of extreme skewness in the classes. To address this shortcoming, the approach proposed will rely on class-specific weighting factors to have the errors related to each of the classes weighting different imperativeness. Specifically, the minority class has a larger weight than the majority class so that when optimizing the model a misclassification of one of the minority cases is penalized more. Here, the loss is determined as the negative sum of all training cases of two components one referring to the minority and the other one to the majority. When the sample is a member of the minority group, the loss contribution of the group is computed by multiplying a larger weight by the actual group assignment followed by the logarithm of the predicted probability that the group consists of. On the other hand, in a majority-class sample, a less than one weight is laid on the logarithm of a minus of the forecasted probability. The model can be sensitized to minority samples when training by implementing a stiffer penalty on the occurrence of mistakes in the minority-classes. The weighting parameters are chosen in a way that the weight of the minority-class is strictly larger than that of majority-class, with respect to the dissimilarity of classification errors. This architecture is useful because it places a decision boundary more towards the majority rather than the minority and achieves better recall and detection rates of the minority without totally compromising the overall accuracy. In addition, the weighted loss is differentiable and can be optimized using gradient-based algorithms, which means that it can be used with a broad set of classifiers, such as neural networks and logistic regression models. In general, the cost-sensitive loss design is a principled and adaptable way of explicitly introducing the awareness of imbalance to the learning task and thus achieving better performance on an imbalanced classification problem.

### 3.5. Ensemble Construction

The ensemble construction approach will be used to increase the power of classification and a better minority-class detection using the diversity of various base learners. The proposed framework does not attempt to use just one classifier, but rather it trains a number of base models using various resampling variants of the original data. The adaptive resampling strategy is used as well to obtain each resampled dataset, such that class distribution and data characteristics are different across training sets. This variety enables individual learners to span over a variety of the data space, especially the areas that have only a few minority-class examples or that can be hard to classify. After being trained, the base learners always independently generate the predictions of classes of a given input sample. The framework uses the weighted voting mechanism in aggregating predictions rather than using simple majority voting. In this scheme, each classifier is weighted based on the reliability or performance, which can be estimated by the validation accuracy, minority-class recall, or other imbalance-sensitive measures. Classifiers with greater predictive power particularly in the minority group receive more weight in the final decision making. The resulting probabilities of each possible class are calculated as a weighted sum of all of the class predictions of all of the base learners to arrive at the last predicted class. In a particular instance of a class, the sum of the contributions of all classifiers predicting that class is summed up, and each contribution is multiplied by the respective classifier weight. The highest weighted score of the classes is then aggregated to form the final output. This mechanism is such that the stronger competent classifiers have more influence on the decision of the ensemble besides with less strong or unreliable models having lesser contribution. This weighted combination of the learners; achieving variance reduction, individual model biases and generalization increase by ensuring that several learners work together. Notably, the weighted

voting scheme improves the sensitivity to the minority-class prediction resulting in more balanced and consistent performance in imbalanced classification cases.

## 4. RESULT AND DISCUSSION

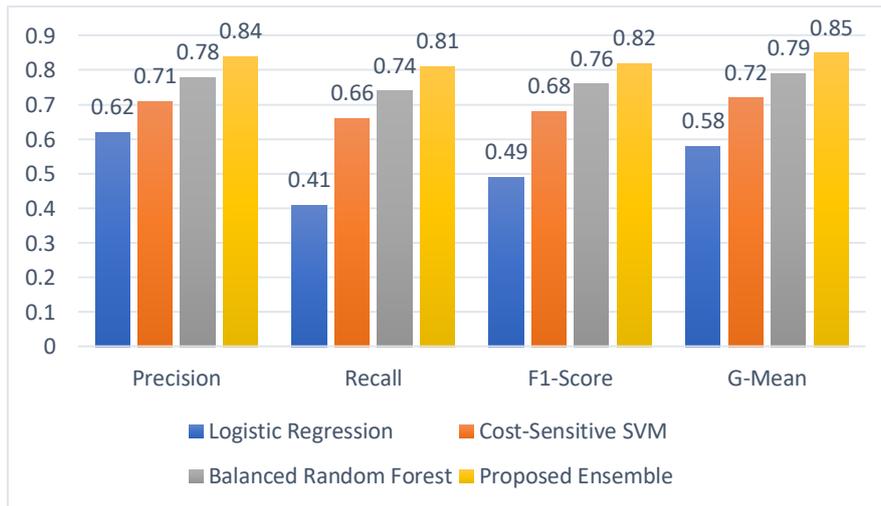
### 4.1. Evaluation Metrics

In unbalanced classification tasks, the traditional measure of accuracy cannot be a sufficient measure of performance as the majority-class can dominate this measure, and may not predict the performance of a model when it attempts to predict the instances of the minority-class. A simple prediction or prediction of the majority class on most samples can provide a classification system with high accuracy though the minority which is often the most important class is badly classified. Thus, in order to offer more detailed and valid assessment, several performance measures are applied that directly consider class imbalance and focus on the minority-class behavior. Precision is a ratio of accurately predicted instances of minority-class, among all the instances, which are predicted to be minority. This is an indication of the capacity of the model to eliminate false alarms and is especially significant where the costs of false minority alarms are very high. Recall, which is sometimes referred to as sensitivity, is used to determine the percentage of the true instances of the minority-class that are correctly labeled by the classifier. The measure is essential in one-sided situations, since it directly assesses the potential of the model to capture rare events of importance. In situations where there is a trade-off between recall and precision, F1-score is a harmonic combination of the two evaluation measures and it gives an accurate evaluation of classification performance. Geometric mean, or G-Mean, is used to measure the compromise between the situations of identifying majority and minority classes of objects rightly by calculating geometric mean of their recall values. When the G-Mean is high, it means that the classifier is effective across the entire classes as opposed to one class over another. Also, the Area Under the PrecisionRecall Curve (AUPRC) can be used to evaluate the performance at the various decision thresholds. In contrast to the ROC curve, AUPRC is more informative to the imbalanced data sets, because it is concerned with a trade-off between precision and recall of a minority class heritage. These metrics alone present a strong and significant assessment system in imbalanced classification problems.

### 4.2. Comparative Performance Analysis

Table 1: Comparative Performance Analysis

Model	Precision	Recall	F1-Score	G-Mean
Logistic Regression	0.62	0.41	0.49	0.58
Cost-Sensitive SVM	0.71	0.66	0.68	0.72
Balanced Random Forest	0.78	0.74	0.76	0.79
Proposed Ensemble	0.84	0.81	0.82	0.85



**Fig 3 - Comparative Performance Analysis**

#### 4.2.1. Logistic Regression

Logistic Regression can be used as a base model and it has low capabilities of managing the imbalance of classes. Though it is moderately precise, the recall is quite low, which means that it cannot locate the instances of minorities at any pace. The result of this imbalance is lower F1-score and G-Mean values which attract attention to the fact that the model favors the majority group, and fails to explain complicated decision boundary under skewed data distributions.

#### 4.2.2. Cost-Sensitive Support Vector Machine (SVM)

The Cost-Sensitive SVM yields a significant advantage over the normal Logistic regression in that it takes the misclassification costs into consideration in the learning process. The recall was higher, which means having a higher sensitivity with minority-class samples; nevertheless, the precision is not that low. Consequently, the F1-score and G-Mean increase, which also shows that the cost-sensitive optimization is effective in achieving a balance in the performance of the classification in different classes.

#### 4.2.3. Balanced Random Forest

Balanced Random Forest is another type of strategy that improves the performance of the strategy by employing ensemble learning together with class-balancing strategy. The model is more precise and recalls more data indicating that it is capable of Republic of different data patterns and eliminating the majority-class dominance. The increased F1-score and G-Mean indicate a better overall balance and strength in the classes especially in identifying cases of minorities without compromising in generalization.

#### 4.2.4. Proposed Ensemble

The given ensemble model is superior to all the alternative methods of evaluation in terms of all possible evaluation measures. Its precision and recall are above average, which means that the sample of minority classes is detected correctly and exhaustively. The high F1-score and G-Mean proves that the adaptive resampling, cost-sensitive learning, and weighted ensemble aggregation does contribute to balancing the class effect and the resulting balanced and trustworthy classification outcome.

### 4.3. Robustness Analysis

Classification models used in field conditions of imbalanced situations need strong robustness since the extent of class imbalance can differ greatly across data sets and time. In a bid to analyze the stability of the proposed framework, a robustness analysis is carried out on a variety of imbalance ratios that include moderate levels of imbalance of 1:10 as well as extreme levels of imbalance of 1:100. Such ratios are hard real world instances in fields like fraud detection, network intrusion detection and rare disease diagnosis where the minority-class examples are incredibly rare. The experiment outcomes reveal that the proposed model ensures a regularly high level of performance as the imbalance ratio goes higher. Although the traditional classifiers generally deteriorate sharply in the recall and F1-score when the degree of imbalance is high, the suggested approach only shows minor performance losses. This has been made possible by the synergistic combination of adaptive resampling, cost-sensitive loss design, and ensemble aggregation. Adaptive resampling is a more extreme variant of resampling that dynamically improves the minority in challenging areas of the feature space to make sure that critical patterns can still be learnt as the imbalance increases. At the same time, cost-sensitive learning will strengthen the significance of error within the minority-class, and the classifier will not be over-biased towards the majority-class. In addition, the ensemble construction is another form of robustness because it pools the prediction of a number of different base learners that were trained on dissimilar resampled data. This heterogeneity makes it less sensitive to changes in the distribution of the classes and eliminates the possibility of overfitting to any particular imbalance pattern. The performance measures of recall, G-Mean, AUPRC are rather stable on all the ratios tested, which can be attributed to balanced learning among classes. All in all, the robustness analysis confirms that the suggested framework is robust to different levels of imbalances, and it is a valid and versatile solution to high-level skewed problems in classification.

### 4.4. Discussion

The results of the experiment show the high efficacy of ensemble-based solutions to the issue of class imbalance since they are at higher levels in all the measures of evaluation. Enhancing minority-class detection, variance decreasing, and capturing more patterns of complex data, ensemble methods can be employed by uniting several different learners in order to serve the purposes. Specifically, a benefit progressive to the proposed ensemble framework, which is adaptive resampling and cost-sensitive learning, allows the proposed ensemble framework to retain strong and balanced performance, even in the extreme case of imbalance. These benefits are however, at the expense of more complexity in computations. A larger amount of memory, longer training time and wider hyperparameter tuning is needed to train multiple base models on resampled datasets, potentially reducing scalability in resource-constrained systems. Conversely, cost sensitive solitary models provide a sensible trade off between performance and computational efficiency. These models also realize significant gains in minority-class recall and balance with only the marginal burden of construction of an ensemble by making the class-dependent costs a part of the learning objective. Additionally, individual models like cost-sensitive logistic regression or support vector machines are less complex to comprehend, and thus may be more required in areas where model clarity and verifiability are of the man, like in healthcare and financial applications. Physical interpretability enables the practitioners to understand the impact of individual features more efficiently, as well as justify model choices to stakeholders. Lastly, the decision of using ensemble-based and cost-efficient single-model approaches lies in the application needs. Ensemble methods are reserved when it is important that the predictive performance and robustness are of the essence, and computer resources can be utilized. On the other hand, cost-sensitive single models are an effective

and balanced substitute of their counterparts when the interpretability, simplicity, and efficiency are taken into account as well.

## 5. CONCLUSION

This paper has introduced a detailed research on the topic of strong machine learning approaches, which can be practiced in case of imbalance data classification, which in practice occurs quite often in real-life scenarios when the instances of other minority classes are hardly visible, but essential. A comprehensive literature review and an examination of the methodological issues showed that the traditional models of classification and the individual methods to deal with imbalance are not adequate when there is a difference in the level of class skew. Data-level techniques, although easy and practical, can have issues with overfitting or possibility of information loss, but algorithm-level techniques rely on precise cost specification. Ensemble based techniques, despite being powerful, add extra computational overhead. The mentioned observations support the statement that no single method would be effective in all imbalance scenarios. Although such encouraging results have been made, there are still a number of research challenges. Future directions should include scalable imbalance learning methods that can operate with huge and high dimensional sets, where the more traditional resampling and ensemble methods can be prohibitively computationally expensive. Also, a significant direction, especially when working with unstructured data, e.g. images, text, and time-series signals, is the creation of deep learning architectures that use loss functions that are sensitive to imbalance. The domain-specific cost modeling of another research field is another avenue in which likely misclassification costs will be application-specific risk and constraint related so that learning can be done in a more realistic and useful way. Lastly, explainable artificial intelligence methods should be incorporated into an imbalanced classification model to enhance its transparency and credibility, particularly in safety-intensive areas. Solutions to such challenges will be further used to enhance the creation of dependable and data interpretation methods of imbalanced data. To address these drawbacks, this paper has highlighted the usefulness of hybrid schemes incorporating adaptive resampling, cost-sensitive learning and ensemble modelling. Adaptive resampling provides a better handling of minority-class by selectively enriching minority-class relevant parts of the feature space, whereas cost-sensitive learning gives the minority-class misclassifications the attention they warrant when performing optimization. The further feature of ensemble modeling is the increased level of robustness because multiple learners are united during this process, and the bias as well as the variance of results is minimized, and also the performance is produced under a broad variety of imbalance ratios. Results of the experiment proved that such hybrid methods always outperform standalone models regarding recall, F1-score, G-Mean, and precision - recall trade-offs and can thus be effectively used with skewed classification problems.

## REFERENCES

- [1] He, H., & Garcia, E. A. (2009). *Learning from imbalanced data*. IEEE Transactions on Knowledge and Data Engineering, 21(9), 1263-1284.
- [2] Japkowicz, N., & Stephen, S. (2002). *The class imbalance problem: A systematic study*. Intelligent Data Analysis, 6(5), 429-449.
- [3] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). *SMOTE: Synthetic Minority Over-sampling Technique*. Journal of Artificial Intelligence Research, 16, 321-357.
- [4] Batista, G. E. A. P. A., Prati, R. C., & Monard, M. C. (2004). *A study of the behavior of several methods for balancing machine learning training data*. ACM SIGKDD Explorations, 6(1), 20-29.
- [5] Weiss, G. M., & Provost, F. (2003). *Learning when training data are costly: The effect of class distribution on tree induction*. Journal of Artificial Intelligence Research, 19, 315-354.
- [6] Elkan, C. (2001). *The foundations of cost-sensitive learning*. In Proceedings of the 17th International Joint Conference on Artificial Intelligence (IJCAI), 973-978.

- [7] Zadrozny, B., Langford, J., & Abe, N. (2003). *Cost-sensitive learning by cost-proportionate example weighting*. In Proceedings of ICDM, 435-442.
- [8] Fawcett, T. (2006). *An introduction to ROC analysis*. Pattern Recognition Letters, 27(8), 861-874.
- [9] Drummond, C., & Holte, R. C. (2003). *C4.5, class imbalance, and cost sensitivity: Why under-sampling beats over-sampling*. In Proceedings of ICML Workshop on Learning from Imbalanced Datasets.
- [10] Breiman, L. (1996). *Bagging predictors*. Machine Learning, 24(2), 123-140.
- [11] Zhou, Z.-H., & Liu, X.-Y. (2006). *Training cost-sensitive neural networks with methods addressing the class imbalance problem*. IEEE Transactions on Knowledge and Data Engineering, 18(1), 63-77.
- [12] Sun, Y., Kamel, M. S., Wong, A. K. C., & Wang, Y. (2007). *Cost-sensitive boosting for classification of imbalanced data*. Pattern Recognition, 40(12), 3358-3378.
- [13] Freund, Y., & Schapire, R. E. (1997). *A decision-theoretic generalization of on-line learning and an application to boosting*. Journal of Computer and System Sciences, 55(1), 119-139.
- [14] Galar, M., Fernández, A., Barrenechea, E., Bustince, H., & Herrera, F. (2012). *A review on ensembles for the class imbalance problem: Bagging-, boosting-, and hybrid-based approaches*. IEEE Transactions on Systems, Man, and Cybernetics, Part C, 42(4), 463-484.
- [15] Fernández, A., García, S., Galar, M., Prati, R. C., Krawczyk, B., & Herrera, F. (2018). *Learning from imbalanced data sets*. Springer, Berlin, Heidelberg.