# Combining Hybrid Approach Redefinition-Multiclass Imbalance (HAR-MI) and Hybrid Sampling in Handling Multi-Class Imbalance and Overlapping

Hartono[a,b,*], Erianto Ongko[c]

[a] Department of Computer Science, Universitas IBBI, Medan, 20114, Indonesia
[b] Department of Computer Science, Universitas Potensi Utama, Medan, 20241, Indonesia
[c] Department of Informatics, Akademi Teknologi Industri Immanuel, 20114, Medan, Indonesia
Corresponding author: *hartono@ibbi.ac.id

*Abstract*—The class imbalance problem in the multi-class dataset is more challenging to manage than the problem in the two classes and this problem is more complicated if accompanied by overlapping. One method that has proven reliable in dealing with this problem is the Hybrid Approach Redefinition-Multiclass Imbalance (HAR-MI) method which is classified as a hybrid approach that combines sampling and classifier ensembles. However, in terms of diversity among classifiers, a hybrid approach that combines sampling and classifier ensembles will give better results. HAR-MI provides excellent results in handling multi-class imbalances. The HAR-MI method uses SMOTE to increase the number of samples in the minority class. However, this SMOTE also has a weakness where an extremely imbalanced dataset and a large number of attributes will be over-fitting. To overcome the problem of over-fitting, the Hybrid Sampling method was proposed. HAR-MI combination with Hybrid Sampling is done to increase the number of samples in the minority class and at the same time reduce the number of noise samples in the majority class. The preprocessing stages at HAR-MI will use the Minimizing Overlapping Selection under Hybrid Sampling (MOSHS) method, and the processing stages will use Different Contribution Sampling. The results obtained will be compared with the results using Neighbourhood-based under-sampling. Overlapping and Classifier Performance will be measured using Augmented R-Value, the Matthews Correlation Coefficient (MCC), Precision, Recall, and F-Value. The results showed that HAR-MI with Hybrid Sampling gave better results in terms of Augmented R-Value, Precision, Recall, and F-Value

*Keywords*— Class imbalance; multi-class dataset; multi-class imbalance; hybrid approach; HAR-MI.

## I. INTRODUCTION

The problem of class imbalance has become one of the most exciting data mining problems [1]. The class imbalance has become one of the most interesting research issues regarding data mining, machine learning, and knowledge discovery[2]. This problem occurs because most of the real-world dataset is in an imbalanced state and if it is not handled properly it will cause a class with a small number of samples to become unrepresented and reduce the level of accuracy[3]. In general, the approach to solving class imbalance problems can be divided into 3 (three), namely: data-level, algorithm-level, and hybrid[4]. The data-level approach focuses on efforts to change the distribution of data through a process of over sampling or under-sampling. Oversampling was carried out on the minority class and under-sampling was carried out

on the majority class[5]. On the other hand, the algorithm-level approach does not change the distribution of data, but focuses on classifier efforts to pay more attention to minority classes by applying bagging, boosting, or through the ensemble process of existing classifiers[6].

Hybrid Approach is an approach that combines Data-Level and Algorithm-Level[7]. In terms of diversity and classifier performance, a hybrid approach that combines sampling and classifier ensembles will give good results[8]. The Hybrid Method is good at dealing with the binary-class imbalance and multi-class imbalance problems[9]. Multi-class imbalance problems are more difficult to handle than binary-class imbalance, and usually, multi-class balance problems do not stand alone but are accompanied by overlapping[10]. This problem becomes even more challenging if the minority classes are in overlapping conditions[11].

To minimize the impact of multi-class imbalance which is accompanied by overlapping, the preprocessing process has a very significant effect[12]. For this problem, the feature selection method is often used at the preprocessing stage, so the effort to apply the preprocessing stage in the hybrid approach is a wise choice[13]. One of the hybrid approach methods that was applied to preprocess and gives satisfactory results in this problem is the Hybrid Approach Redefinition-Multiclass Imbalance (HAR-MI)[14].

As with most hybrid approach methods, HAR-MI also uses the oversampling method for minority classes by using SMOTE in the feature selection process at the preprocessing stage. One of the Feature Selection methods that provide excellent results in handling overlapping is Minimizing Overlapping Selection under SMOTE (MOSS)[15], even though this oversampling process often causes overfitting[16]. Besides, other problems that are often found in the application of SMOTE are overgeneralization and noise[17]. The use of Minority Over-Sampling Techniques (M-SMOTE) and Edited Nearest Neighbor (ENN), which are a type of Hybrid Sampling, has yielded very satisfying results [18].

It would be interesting if there is a method that combines multi-class balance handling followed by overlapping and at the same time paying attention so that the sampling process does not overfit. This study will combine the use of HAR-MI with Hybrid Sampling. This study's results will be compared with Neighborhood-based under-sampling, which is one of the best methods of handling multi-class imbalance and overlapping[19].

## II. Materials and Method

### A. Hybrid Approach

The pseudocode of the Hybrid Approach is as follows[20].

$Input: D_T = \{x_1, x_2, \ldots, x_n\}//Training\ Dataset$
$N = Number\ of\ Classifier$
$Output: Classification\ Prediction\ P$
$Method:$
$Step\ 1\ Preprocessing\ using\ Preprocessing\ Method$
$Step\ 2\ For\ i = 1\ to\ N\ do$
$\quad i. Apply\ Machine\ Learning\ Classification\ Algorithm$
$\quad on\ The\ Attributes\ of\ D_T$
$\quad ii. Obtain\ Classification\ Prediction\ P_i\ from\ machine$
$\quad learning\ classification\ algorithm$
$End\ For$
$Step\ 3\ For\ i = 1\ to\ n$
$\quad Apply\ processing\ using\ bagging, boosting\ or\ sampling$
$End\ For$

### B. Hybrid Sampling

The pseudocode of the Hybrid Sampling using M-SMOTE and ENN is as follows[18].

$Input: Dataset\ S, Minority\ Samples\ S_{Min}, Majority\ Sample\ S_{Maj}$
$Output: Final\ Dataset\ S'$
$Create\ global\ variable\ G_{max}, create\ array\ Eva_{min}, Eva_{maj}, Eva$

$Step\ 1: If\ G_{MCC} = 0$
$\quad Processing\ M - SMOTE\ for\ S_{Min}$
$\quad Processing\ ENN\ for\ S_{Maj}$
$\quad End\ If$
$Step\ 2: Calculate\ Eva_{min}\ using\ MCC$
$\quad Calculate\ Eva_{maj}\ using\ MCC$
$\quad Calculate\ Eva\ using\ MCC$
$Step\ 3: If\ Eva_{min} < Eva\ or\ If\ Eva_{maj} < Eva$
$\quad G_{MCC} = G_{MCC} - 1$
$\quad End\ If$
$Step\ 4: if\ G_{MCC} < 0$
$\quad Terminate\ and\ Output\ Final\ Dataset\ S'$
$\quad else$
$\quad Return\ to\ Step\ 1$
$\quad End\ If$

### C. Augmented R-Value

Augmented R-Value states how much overlapping occurs. The greater the Augmented R-Value, the greater the overlapping[21].

$$R_{aug}(D[V]) = \frac{\sum_{i=0}^{k-1}|C_{k-1-i}|R(C_i)}{\sum_{i=0}^{k-1}|c_i|} \quad (1)$$

Where $C_0, C_1, \ldots, C_{k-1}$ are $k$ class labels with $|C_0| \geq |C_1| \geq \cdots \geq |C_{k-1}|$ and $D[V]$: Dataset D containing predictors in set $V$. Larger $R_{Aug}$ is higher overlap degree of a dataset.

### D. Classifier Performance

Classifier Performance was measured using the Matthews Correlation Coefficient (MCC), Precision, Recall, and F-Value. This classifier performance measurement is carried out based on the confusion matrix shown in Table 1[22].

TABLE I
CONFUSION MATRIX

|  |  | Predictive Positive Class | Predictive Negative Class |
|---|---|---|---|
| Actual Class | Positive | True Positive (TP) | False Negative (FN) |
| Actual Class | Negative | False Positive (FP) | True Negative (TN) |

The Matthews Correlation Coefficient (MCC), Precision, Recall, and F-Value calculations can be seen in the following equation[18].

$$MCC = \frac{TP\ x\ TN - FP\ x\ FN}{\sqrt{(TN\ x\ FN)(TN\ x\ FP)(TN\ x\ FN)(TP\ x\ FP)}} \quad (2)$$

$$Precision = \frac{TP}{TP + FP} \quad (3)$$

$$Recall = TP \quad (4)$$

$$F - Value = \frac{2\ x\ Precision\ x\ Recall}{Precision + Recall} \quad (5)$$

### E. Proposed Method / Algorithm

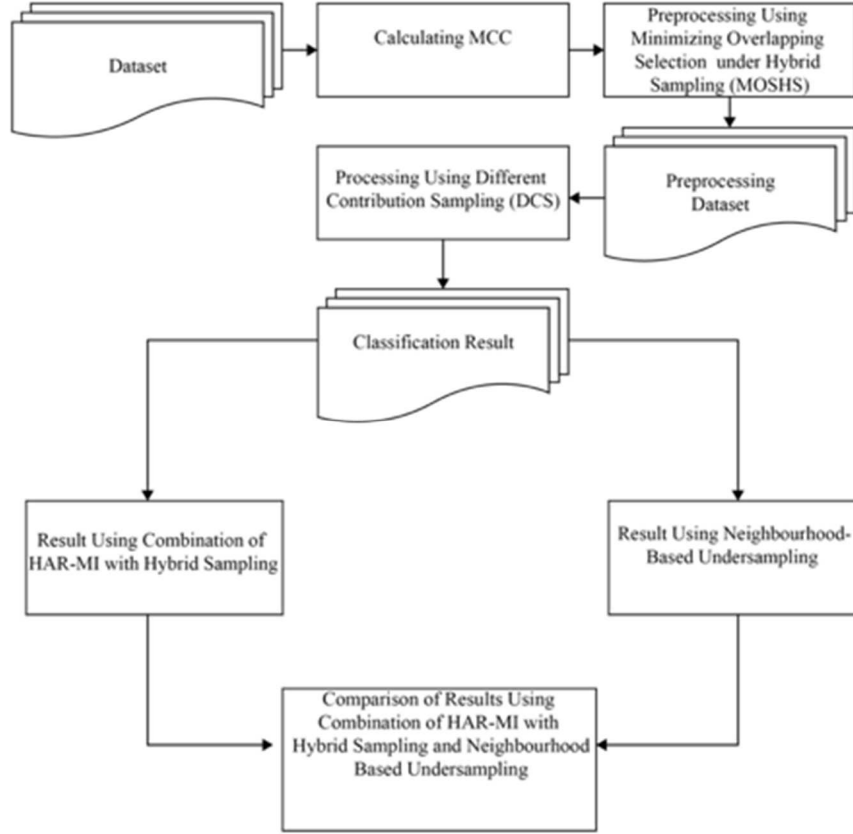The research stages can be seen in Fig. 1.

23

Fig. 1  Research Stage

## F. Preprocessing Using Minimizing Overlapping Selection under Hybrid Sampling (MOSHS)

The pseudocode of the preprocessing stage is as follows.

1: $X - matrix\ with\ p\ predictors: X = [x_1, x_2, ..., x_p]; class\ label: y$
2: $For\ All\ Samples\ in\ Minority$
3: $\quad Hybrid\ Sampling\ the\ Minority\ Class\ using\ m - SMOTE$
4: $End\ For$
5: $Create\ NewMinority$
6: $For\ All\ Samples\ in\ Majority$
7: $\quad Hybrid\ Sampling\ the\ Majority\ Class\ using\ ENN$
8: $End\ For$
9: $Create\ NewMajority$
10: $ForAll\ Samples\ in\ NewMinority\ and\ NewMajority$
11: $\quad Preprocessed\ Dataset$
12: $End\ For$

## G. Processing Using Different Contribution Sampling (DCS)

The pseudocode of the processing stage is as follows.

1: $For\ i = 1\ to\ Number\ of\ Instance\ in\ Preprocessed\ Dataset$
2: $\quad Add\ Preprocessed\ Dataset\ to\ S_i$
3: $\quad B - SVM\ will\ do\ for\ Classifying\ S_i$
4: $\quad Determine\ the\ Majority\ Class$
5: $\quad Determine\ the\ Minority\ Class$
6: $\quad For\ All\ Instance\ in\ Majority\ Class$
7: $\qquad NewSVSets[\ ]\ will\ form\ by\ checking\ and\ delete\ the\ noise\ in\ SVSets$
8: $\qquad NewNSVSets[\ ]\ will\ form\ by\ Multiple\ Hybrid\ Sampling$
9: $\quad End\ For$
10: $\quad For\ All\ Instance\ from\ NewSVSets\ and\ NSVSets$
11: $\qquad Create\ an\ instance\ for\ majority\ class$
12: $\quad End\ For$
13: $\quad For\ All\ Instance\ in\ minority\ class$
14: $\qquad SMOTEBoost\ Process\ for\ SVSets\ and\ Create\ SMOTESets$
15: $\quad End\ For$
16: $\quad For\ All\ SMOTESets\ and\ NewNSVSets\ do$
17: $\qquad NewPositiveSampleSets$
18: $\quad End\ For$
19: $\quad For\ All\ NewNegativeSampleSets\ and\ NewPositiveSampleSets\ do$
20: $\qquad ResultDataSet$
21: $\quad End\ For$
22: $End\ For$

## III. RESULTS AND DISCUSSION

### A. Dataset Description

The multi-class imbalanced datasets used in this study were sourced from the KEEL Repository[23]. The dataset used can be seen in Table II.

TABLE II
DATASET DESCRIPTION

| Dataset | #Ex | #Atts | Distribution of Class | IR |
|---|---|---|---|---|
| Contraceptive | 1473 | 9 | 629/333/511 | 1.89 |
| Flare | 1066 | 11 | 147/211/239/95/43/331 | 7.70 |
| Car Evaluation | 1728 | 6 | 384/69/1210/65 | 18.62 |
| Thyroid Disease | 720 | 21 | 17/37/666 | 39.18 |
| Red Wine Quality | 1599 | 11 | 10/53/681/638/199/18 | 68.10 |
| Page-Blocks | 5473 | 10 | 4913/329/28/88/ 115 | 188.72 |

Table II shows that the dataset used has various imbalance ratios, ranging from low, medium, and high imbalance ratios. Likewise, the number of samples also varied.

24

## B. Testing Result

The first test was conducted to obtain Augmented R-Value and MCC values. The test results can be seen in Table III.

TABLE III
TESTING FOR AUGMENTED R-VALUE AND MCC

| Dataset | HAR-MI with Hybrid Sampling | | Neighborhood Based Under-sampling | |
| --- | --- | --- | --- | --- |
| | Augmented R-Value | MCC | Augmented R-Value | MCC |
| Contraceptive | 0.327 | 0.97 | 0.337 | 0.91 |
| Flare | 0.357 | 0.83 | 0.359 | 0.82 |
| Car Evaluation | 0.367 | 0.85 | 0.373 | 0.81 |
| Thyroid Disease | 0.379 | 0.81 | 0.381 | 0.79 |
| Red Wine Quality | 0.411 | 0.75 | 0.415 | 0.71 |
| Page-Blocks | 0.436 | 0.73 | 0.437 | 0.71 |

Based on Table III, it can be seen that for the Augmented R-Value, the results obtained by HAR-MI with Hybrid Sampling are better than the Neighborhood-based under-sampling. The greater the Augmented R-Value, the greater the overlapping that occurs. Based on the Augmented R-Value obtained by the two methods, the greater the imbalance ratio value, the greater the tendency for overlapping to occur. The MCC value provided by HAR-MI with Hybrid Sampling is also better than that obtained by Neighborhood-based under-sampling. The second test was conducted to obtain Precision, Recall, and F-Value. The test results can be seen in Table IV.

TABLE IV
TESTING FOR PRECISION, RECALL, AND F-VALUE

| Dataset | HAR-MI with Hybrid Sampling | | | Neighborhood Based Under-sampling | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Precision | Recall | F-Value | Precision | Recall | F-Value |
| Contraceptive | 0.88 | 0.97 | 0.92 | 0.78 | 0.89 | 0.83 |
| Flare | 0.85 | 0.88 | 0.87 | 0.81 | 0.87 | 0.84 |
| Car Evaluation | 0.84 | 0.89 | 0.86 | 0.76 | 0.73 | 0.75 |
| Thyroid Disease | 0.87 | 0.76 | 0.81 | 0.85 | 0.71 | 0.77 |
| Red Wine Quality | 0.82 | 0.81 | 0.81 | 0.82 | 0.72 | 0.77 |
| Page-Blocks | 0.78 | 0.77 | 0.77 | 0.77 | 0.69 | 0.73 |

Based on Table IV, it can be seen that based on the Precision, Recall, and F-Value values the results given by HAR-MI with Hybrid Sampling are better than the results obtained by Neighborhood-based under-sampling.

## C. Statistical Tests

To validate the results of the study, a statistical test was conducted to measure performance using the Wilcoxon Signed-Rank Test[24]. The statistical test results can be seen in Table V.

TABLE V
STATISTICAL TESTS USING WILCOXON SIGNED-RANK TEST

| Performance Measurement | P-Value | Hypothesis |
| --- | --- | --- |
| Augmented R-Value | 0.0355223 | $H_0$ (no significant score difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling) rejected and this means $H_1$ (there is a significant difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling in score) Accepted because the p-value <0.05 |
| MCC | 0.0355223 | $H_0$ (no significant score difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling) rejected and this means $H_1$ (there is a significant difference between HAR-MI with Hybrd Sampling and Neighbourhood-Based Under-sampling in score) Accepted because the p-value <0.05 |
| Precision | 0.0625000 | $H_0$ (no significant score difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling) is accepted and this means $H_1$ (there is a significant difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling in score) is rejected because the p-value >0.05 |
| Recall | 0.0312500 | $H_0$ (no significant score difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling) rejected and this means $H_1$ (there is a significant difference between HAR-MI with Hybrd Sampling and Neighbourhood-Based Under-sampling in score) Accepted because the p-value <0.05 |
| F-Value | 0.0340064 | $H_0$ (no significant score difference between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling) rejected and this means $H_1$ (there is a significant difference between HAR-MI with Hybrd Sampling and Neighbourhood-Based Under-sampling in score) Accepted because the p-value <0.05 |

## D. Discussion

Based on the test results and Statistical Tests, it can be seen that in terms of overlapping the HAR-MI method with Hybrid Sampling gives better results compared to MCC between HAR-MI with Hybrid Sampling and Neighborhood-Based Under-sampling. However, in general, the results obtained in overlapping handling are good, where the Augmented R-Value obtained is not too high. Augmented R-Value is very dependent on the imbalance ratio; the higher the value of the imbalance ratio, the higher the overlapping that occurs. There is a significant difference for Augmented R-Value and MCC between HAR-MI with Hybrid Sampling and Neighbourhood-Based Under-sampling based on statistical tests.

As for the MCC value, the results given by HAR-MI with Hybrid Sampling are still better and there is a tendency that the more classes there are, the lower the MCC value obtained. As for the Precision, Recall, and F-Value values, the results obtained show that HAR-MI with Hybrid Sampling is also better than MCC between HAR-MI with Hybrid Sampling

and Neighbourhood-Based Under-sampling. The results obtained show that the higher the imbalance ratio, the value of Precision, Recall, and F-Value obtained also decreases.

Based on the results of statistical testing with the Wilcoxon Signed-Rank Test, it was found that for Augmented R-Value, the P-Value is 0.0355223, the P-Value for MCC is 0.0355223, the P-Value for Recall is 0.0312500, and the P-Value for F-Value is 0.0340064. This means that for Augmented R-Value, MCC, Recall, and F-Value, there is a significant difference between HAR-MI results with Hybrid Sampling and Neighborhood-Based Under-sampling. As for Precision, although HAR-MI results are better than Neighborhood-Based Under-sampling but based on the test results with the Wilcoxon Signed-Rank Test, there is no significant difference as indicated by the P-Value obtained> 0.05, where the P-Value obtained is 0.0625000.

## IV. CONCLUSION

Based on the results in Tables III, IV, and V, it can be seen that in terms of handling multi-class imbalance and overlapping, the results obtained using HAR-MI with Hybrid Sampling give better results compared to Neighbourhood-Based Under-sampling. The results obtained show that HAR-MI with Hybrid Sampling excels at all test values such as Augmented R-Value, MCC, Precision, Recall, and F-Value.

This shows that for handling multi-class imbalance, Hybrid Sampling, which can avoid over fitting, also gives better results compared to Under-sampling or Over Sampling. Future Research can pay attention to the handling of multi-class imbalance accompanied by overlapping in a state of high yield ratio and datasets with a large number of classes and many attributes.

## REFERENCES

[1] G. Haixiang, L. Yijing, J. Shang, G. Mingyun, H. Yuanyue, and G. Bing, "Learning from Class-Imbalanced Data: Review of Methods and Applications," *Expert Systems with Applications*, vol. 73, pp. 220–239, May 2017.

[2] A. Guzmán-Ponce, J. S. Sánchez, R. M. Valdovinos, and J. R. Marcial-Romero, "DBIG-US: A two-stage under-sampling algorithm to face the class imbalance problem," *Expert Systems with Applications*, p. 114301, Nov. 2020, doi: 10.1016/j.eswa.2020.114301.

[3] B. Liu and G. Tsoumakas, "Dealing with class imbalance in classifier chains via random under-sampling," *Knowledge-Based Systems*, vol. 192, p. 105292, Mar. 2020, doi: 10.1016/j.knosys.2019.105292.

[4] J. M. Johnson and T. M. Khoshgoftaar, "Survey on deep learning with class imbalance," *J Big Data*, vol. 6, no. 1, p. 27, Mar. 2019, doi: 10.1186/s40537-019-0192-5.

[5] P. Shamsolmoali, M. Zareapoor, L. Shen, A. H. Sadka, and J. Yang, "Imbalanced data learning by minority class augmentation using capsule adversarial networks," *Neurocomputing*, Jul. 2020, doi: 10.1016/j.neucom.2020.01.119.

[6] W. Hou, X. Wang, H. Zhang, J. Wang, and L. Li, "A novel dynamic ensemble selection classifier for an imbalanced data set: An application for credit risk assessment," *Knowledge-Based Systems*, vol. 208, p. 106462, Nov. 2020, doi: 10.1016/j.knosys.2020.106462.

[7] F. Rayhan, S. Ahmed, A. Mahbub, R. Jani, S. Shatabda, and D. M. Farid, "CUSBoost: Cluster-Based Under-Sampling with Boosting for

Imbalanced Classification," in *2017 2nd International Conference on Computational Systems and Information Technology for Sustainable Solution (CSITSS)*, Dec. 2017, pp. 1–5, doi: 10.1109/CSITSS.2017.8447534.

[8] J. Zhao, J. Jin, S. Chen, R. Zhang, B. Yu, and Q. Liu, "A weighted hybrid ensemble method for classifying imbalanced data," *Knowledge-Based Systems*, vol. 203, p. 106087, Sep. 2020, doi: 10.1016/j.knosys.2020.106087.

[9] Z. Liu, D. Tang, Y. Cai, R. Wang, and F. Chen, "A hybrid method based on ensemble WELM for handling multi class imbalance in cancer microarray data," *Neurocomputing*, vol. 266, pp. 641–650, Nov. 2017, doi: 10.1016/j.neucom.2017.05.066.

[10] E. R. Q. Fernandes and A. C. P. L. F. de Carvalho, "Evolutionary inversion of class distribution in overlapping areas for multi-class imbalanced learning," *Information Sciences*, vol. 494, pp. 141–154, Aug. 2019, doi: 10.1016/j.ins.2019.04.052.

[11] Y. Zhu, Y. Yan, Y. Zhang, and Y. Zhang, "EHSO: Evolutionary Hybrid Sampling in overlapping scenarios for imbalanced learning," *Neurocomputing*, vol. 417, pp. 333–346, Dec. 2020, doi: 10.1016/j.neucom.2020.08.060.

[12] P. Zyblewski, R. Sabourin, and M. Woźniak, "Preprocessed dynamic classifier ensemble selection for highly imbalanced drifted data streams," *Information Fusion*, vol. 66, pp. 138–154, Feb. 2021, doi: 10.1016/j.inffus.2020.09.004.

[13] L. Yijing, G. Haixiang, L. Xiao, L. Yanan, and L. Jinling, "Adapted ensemble classification algorithm based on multiple classifier system and feature selection for classifying multi-class imbalanced data," *Knowledge-Based Systems*, vol. 94, pp. 88–104, Feb. 2016, doi: 10.1016/j.knosys.2015.11.013.

[14] H. Hartono, Y. Risyani, E. Ongko, and D. Abdullah, "HAR-MI method for multi-class imbalanced datasets," *TELKOMNIKA (Telecommunication Computing Electronics and Control)*, vol. 18, no. 2, Art. no. 2, Apr. 2020, doi: 10.12928/telkomnika.v18i2.14818.

[15] G.-H. Fu, Y.-J. Wu, M.-J. Zong, and L.-Z. Yi, "Feature selection and classification by minimizing overlap degree for class-imbalanced data in metabolomics," *Chemometrics and Intelligent Laboratory Systems*, vol. 196, p. 103906, Jan. 2020, doi: 10.1016/j.chemolab.2019.103906.

[16] X. Gao *et al.*, "An ensemble imbalanced classification method based on model dynamic selection driven by data partition hybrid sampling," *Expert Systems with Applications*, vol. 160, p. 113660, Dec. 2020, doi: 10.1016/j.eswa.2020.113660.

[17] J. Wei, H. Huang, L. Yao, Y. Hu, Q. Fan, and D. Huang, "New imbalanced bearing fault diagnosis method based on Sample-characteristic Oversampling TechniquE (SCOTE) and multi-class LS-SVM," *Applied Soft Computing*, vol. 101, p. 107043, Mar. 2021, doi: 10.1016/j.asoc.2020.107043.

[18] Z. Xu, D. Shen, T. Nie, and Y. Kou, "A hybrid sampling algorithm combining M-SMOTE and ENN based on Random forest for medical imbalanced data," *Journal of Biomedical Informatics*, vol. 107, p. 103465, Jul. 2020, doi: 10.1016/j.jbi.2020.103465.

[19] P. Vuttipittayamongkol and E. Elyan, "Neighbourhood-based under-sampling approach for handling imbalanced and overlapped data," *Information Sciences*, vol. 509, pp. 47–70, Jan. 2020, doi: 10.1016/j.ins.2019.08.062.

[20] M. Galar, A. Fernandez, E. Barrenechea, H. Bustince, and F. Herrera, "A Review on Ensembles for the Class Imbalance Problem: Bagging-, Boosting-, and Hybrid-Based Approaches," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 42, no. 4, pp. 463–484, Jul. 2012, doi: 10.1109/TSMCC.2011.2161285.

[21] S. Oh, "A new dataset evaluation method based on category overlap," *Comput. Biol. Med.*, vol. 41, no. 2, pp. 115–122, Feb. 2011, doi: 10.1016/j.compbiomed.2010.12.006.

[22] A. Luque, A. Carrasco, A. Martín, and A. de las Heras, "The impact of class imbalance in classification performance metrics based on the binary confusion matrix," *Pattern Recognition*, vol. 91, pp. 216–231, Jul. 2019, doi: 10.1016/j.patcog.2019.02.023.

[23] J. Alcalá-Fdez *et al.*, "KEEL: a software tool to assess evolutionary algorithms for data mining problems," *Soft Comput*, vol. 13, no. 3, pp. 307–318, Feb. 2009, doi: 10.1007/s00500-008-0323-y.

[24] F. Wilcoxon, "Individual Comparisons by Ranking Methods on JSTOR," *Biometrics Bulletin*, vol. 1, no. 6, pp. 80–83, 1945.