# Using ensemble neural network based on sampling for multiclass classification

*Bain Khusnul* Khotimah[1*]*, Eko* Setiawan[2], *Devie Rosa* Anamisa[1], *Oktavia Rahayu* Puspitarini[3]

[1]Faculty of Engineering, Universitas Trunojoyo Madura, East Java, Indonesia, 69162
[2]Faculty of Agriculture, University of Trunojoyo Madura, East Java, Indonesia, 69162
[3]Faculty of Animal Husbandry, Islamic University of Malang, East Java, Indonesia, 65144

**Abstract.** Multiclass data classification with class imbalance causes classification performance to decrease, especially in the Neural network method. Research shows that the model proposed by eNN can improve model performance for imbalanced data in the selection of superior quality in beef and cattle data. The results of the Ensemble ANN study with adaboost are able to understand complex relationships by measuring the level of correlation with the target class produced. This study aims to overcome the problem of data imbalance in the ensemble neural network method by comparing the oversampling method with undersampling, so that more representative synthetic data is obtained. Performance evaluation is processed using precision, recall and accuracy calculations. Research on superior local Madura cattle data The RUS-eNN method produces the highest average accuracy value compared to others, reaching 98.00% with a recall value of 100%. While the ROS-eNN method produces a difference in accuracy value that is not so far away, namely 97.69%. The research on the sampling-based eNN approach has better accuracy than without using data replication in improving its performance.

## 1 Introduction

The class imbalance problem has become a common problem due to its classification difficulty caused by the imbalanced class distribution. In particular, many ensemble methods have been proposed to deal with the imbalance. However, most of the efforts have focused on the two-class imbalance problem. There are unsolved problems in the multiclass imbalance problem in real-world applications. An accurate collection of diverse neural networks provides better results and fewer errors than a single neural network. Diversity can be achieved by manipulating input data or output data [1]. NN uses ensemble methods to fine-tune model performance and measure a model's uncertainty. The ensemble method with the ANN model can produce ensemble members by varying the initial weights, ANN topology or architecture, training algorithm, and training data [2]. Examples of ensemble diversity using input data are bagging [3] and boosting [4] neural networks. Bagging provides diversity by randomly resampling the original training data into multiple training sets while boosting provides diversity by manipulating each set to output multiple output data. This paper studies the challenges of the multiclass imbalance problem and investigates the generalization ability of several ensemble solutions, including our recently proposed algorithm AdaBoost.NC, with the aim of effectively and directly handling multiclass imbalance. AdaBoost.NC is applied to several real-world multiclass imbalance tasks and

compared with other popular ensemble methods. AdaBoost.NC is better at recognizing examples of minority classes, which can balance the performance between classes. Multiclass neural network classification involves the construction of a neural network that maps input feature vectors to network outputs that typically contain more than two classes [5, 6].

In general, neural network architectures are used to classify multiclass. The approach is generally built by utilizing multiple binary neural networks where each network can be modeled independently. One of the advantages of using this technique is that different features can be applied to train different neural networks [6]. However, each neural network is trained only based on local knowledge, which can result in overlaps or gaps in the classification boundary zone [7]. Another approach is the implementation of a single neural network with multiple outputs. The complexity of this approach is usually high [8, 9]. However, the classification boundary is sharp, which is used to avoid uncertainty in the classification boundary zone [10].

In previous studies, the ANN method has been widely applied in classification, particularly in diagnosing plant diseases. For instance, it has been used to diagnose corn diseases with an impressive accuracy rate of 95.08% [11]. The concept of ensembles and the bias/variance trade-off of ensemble predictions in environmental modeling has shown that the ensemble NN model is more powerful than the single NN model

---

* Corresponding author: bain@trunojoyo.ac.id

and improves generalization NN model capabilities [12]. NN models contain large degrees of freedom and several uncertainty sources influencing output. The uncertainties arising from the input vector and the structure of the NN network are investigated by Asefa [13], which used a sampling and ensemble-based NN approach to overcome data imbalance and uncertainty in ANN output arising from variations in initial parameters and training data [14]. This paper implements a diverse neural network ensemble using multiclass neural networks trained with the same input feature vector but disagreeing on the target codeword. The first network predicts the accurate membership degree, and the second predicts the false one [15]. The boundary between these two predicted outputs may not be sharp. Uncertainties may occur in the boundary zone. This paper also estimates these uncertainties and represents them regarding uncertainty memberships. These three memberships form a neutrosophic set of intervals [16]. The final classification decision is decided from these three memberships. One of the systems with artificial intelligence that has been successfully developed is the machine learning process (computer). A system with artificial intelligence can experience machine learning that imitates the workings of the human biological nervous system, an artificial neural network system (ANN) [17]. Diagnosing onion diseases for farmers has a vital role, not only in a high level of accuracy but also in a high level of speed in diagnosing plant diseases. ANN is an artificial intelligence system that can analyze complex data and has an undefined or non-linear relationship between variables [18, 19]. The process of selecting superior cattle seeds is based on the performance of Madura cattle body condition characteristics as a selection criterion [4]. The quality of cattle seeds can be assessed from various aspects, including body size, cattle health from all types of diseases, and physical defects of cattle in accordance with Indonesian national standards in good beef cattle breeding guidelines [20].

This research focuses on applying the NN ensemble model in identifying onion diseases with imbalanced data and many features, and the hidden layer structure has produced good performance. The performance of NN in an ensemble depends on several factors, such as (1) The dataset used in this study was obtained from observations for multiclass classification. (2) To validate the proposed method's effectiveness, two experimental scenarios were carried out: first, the ANN algorithm was directly used for modeling without considering class imbalance. (3) The last scenario, using oversampling and undersampling, was used to increase the number of datasets to achieve a balanced dataset. Meanwhile, performance measurement in the system is based on calculating the F-measure, recall, precision, and accuracy values. This research will combine data improvement with sampling techniques and appropriate computational improvements to increase accuracy.

## 2 Research methodology

### 2.1 Biology multiclass classification

Classification of living things for various purposes in various fields, namely agriculture, animal husbandry, health, and etc. The problem that arises in the fields of agriculture and animal husbandry is how to determine superior quality, using classification [11]. Determination of superior seeds is very urgent, for example red onion data. Red onions are one type of commodity that is very much needed by the community. However, red onions are very easy to experience changes in quality such as weight loss, changes in volatility and damage because they have a high water content, so good storage methods are needed. Determination of the quality of superior onions is carried out based on their characteristics, level of dryness, no traces of pests, maximum weight, shape and level of dryness [21, 22].

Furthermore, the research on cases of determining the quality of living things such as the case of determining the superior quality of Madura cattle [23]. Cattle are an important source of animal food globally, and each country has unique endemic cattle breeds. However, categorizing cattle, especially in countries such as Indonesia with large cattle populations, is challenging due to the cost and subjectivity of using human experts [24]. This study uses cattle trait data to determine the quality of superior cattle, in order to address the need to maintain the quality of local cattle. This study is to develop a model that can accurately detect and classify superior local Madurese cattle breeds, such as Sonok cattle, Karapan cattle and beef cattle [25].

### 2.2 Inbalanced class

Most of the conventional models assign the majority class to the data and ignore the minority class because of the skewness of the data [7]. One method used to overcome the imbalance class problem is sampling. The sampling method modifies the distribution of data between the majority and minority classes in the training dataset to balance the amount of data for each class. At the data preprocessing stage, it was identified that the dataset used in this study has a considerable imbalance class problem where data instances with large rating values ($>5$) are far fewer than data instances with small/medium rating values. So, we need a preprocessing method to overcome the imbalance class problem [8-9]. Biological data often experiences imbalanced cases, resulting in issues with model inaccuracy that is less than optimal for all data classes [10]. Several solutions that can handle imbalanced cases include oversampling methods [11, 12], undersampling [13, 14], and a hybrid of oversampling and undersampling methods. Oversampling is done by making a replica (resample) of minority data. In contrast, the undersampling method reduces majority data so that more balanced majority and minority data are obtained [15]. Excessive oversampling methods can cause overfitting, while excessive undersampling can

result in losing some vital information in the dataset [16].

The Synthetic Minority Oversampling Technique (SMOTE) is a well-known oversampling method that effectively handles synthetic data overfitting through the K-NN approach with the use of variables [17, 20]. To further enhance the performance of the oversampling method, some researchers have incorporated the undersampling method as a cleaning method [19]. This combination is expected to produce cleaner processed data that is free from noise, thereby boosting the ability of the oversampling method to create synthetic data without replicating noisy data. The undersampling methods used in this research, such as the Neighborhood Cleaning Rule (NCL) [20, 21], play a crucial role in this process.

### 2.2.1 RUS (Random Under-Sampling)

The algorithm overcomes class imbalance by balancing the data distribution with Oversampling and Undersampling algorithm techniques. The working principle is an algorithm-level approach, namely modifying existing algorithms to consider the meaning of minor classes or developing new algorithms. Combine algorithmic approaches and data-level approaches. Undersampling generates random subsamples from majority class instances by selecting samples in the majority class and adding them to the minority class, forming a new training dataset.
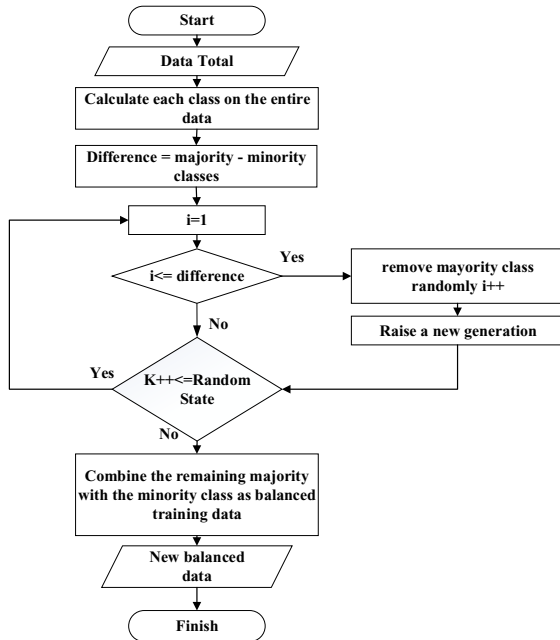


**Fig. 1.** The process of balancing data using Under-Sampling.

### 2.2.2 Random Over Sampling (ROS)

ROS is an over-sampling method in which data in the minority class is reproduced using synthetic data derived from data replication in the minority class. Over-sampling takes an instance of the minority class,

looks for the k-nearest neighbor of each instance, and then generates a synthetic instance instead of replicating the minority class instance. Therefore, it can avoid the problem of excessive overfitting [19]. The algorithm will take the difference value between the vector of features in the minority class and the nearest neighbor value of the minority class and then multiply that value by a random number between 0 and 1. Furthermore, the calculation results are added to the feature vector to obtain the new vector value results [20, 24].

$$X_{new} = X_i + \left(\widehat{X_i} - X_i\right)x\delta \tag{1}$$

with
$X_i$ = vector of features in minority class
$\widehat{X_i}$ = k-nearest neighbors for $Xi$
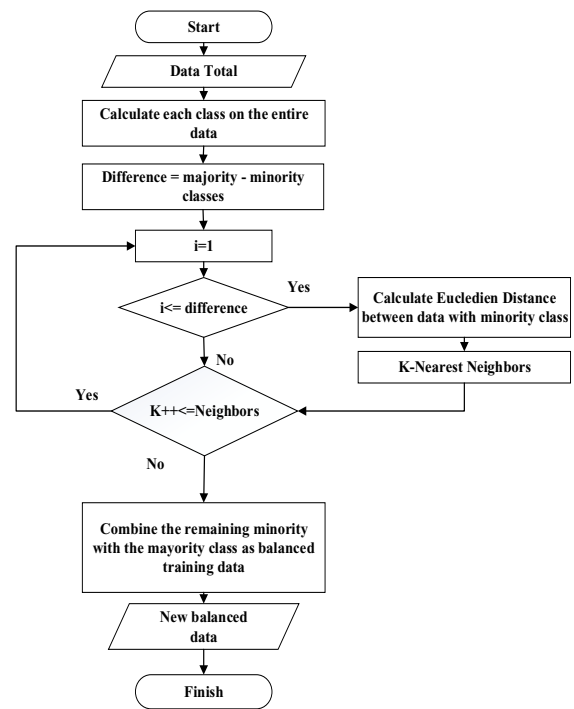$\delta$ = random number between 0 to 1.



**Fig. 2.** The process of balancing data using Over-Sampling.

### 2.3 Ensemble learning

The approach assigns different weights to each method using the least squares method, which optimizes the contribution of each particular individual estimate [25]. Several ensemble learning methods are widely used, such as Boosting, Bagging, and random forest [26]. Boosting is an approach to machine learning to increase accurate predictions by combining many weak and inaccurate rules. Adaptive boosting (AdaBoost) is one of several variants of the boosting algorithm, which is generally combined with a classifier to improve classification performance. The AdaBoost algorithm corresponds to the following steps:
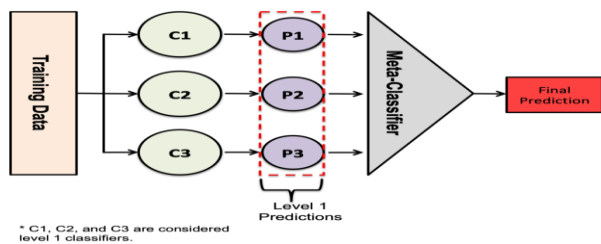
**Fig. 3.** Boots Ensemble learning using of voting technique.

Process:
1. Divide the information into sets for training and validation.
2. The training set should be divided into K folds, for instance 10.
3. Make predictions on the tenth fold after training a base model, such as an SVM, on nine folds.
4. Continue doing this until you have a forecast for every fold.
5. Fit all of the training set with the base model. Utilize the model to forecast the test set.
6. For alternative base models (such as decision trees), repeat steps 3 through 6.
7. Utilize the test set's predictions as features for the meta-model, a new model.
8. Utilizing the meta model, make final predictions for the test set.

## 2.4 Neural network

The most popular neural network algorithm is backpropagation; the algorithm learns on a multi-layer feed-forward neural network consisting of three layers: the input, hidden, and output layers [17]. In some Neural network diagrams, it is possible to have more than one hidden layer, although most contain only one, which is considered sufficient for various purposes [18]. The neural network algorithm for classification uses momentum, a set of weights that can model the data to minimize the average squared distance between the predicted network class and the actual class label of the data tuples [19]. Each training data observation is processed through the network; the input node generates the output value. This output value is then compared with the actual value of the target variable for the resulting error [20].
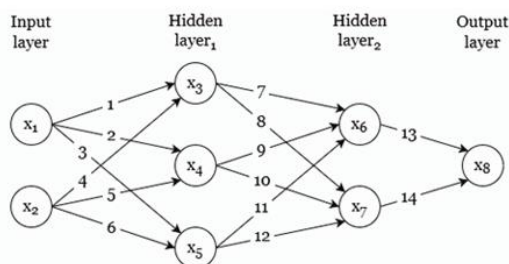


**Fig. 4.** Neural multilayer network.

Fig. 4 explain networks with many layers have hidden layers and weights between the input and output layers. Multilayer Neural Networks can solve more complex problems than networks with a single layer.

## 3 Result and discussion

### 3.1 Preparation of dataset

The data used in this study are onion and cattle data, which are included in the multiclass inbalance to select superior quality in agriculture and animal husbandry. The data were taken from the agriculture and animal husbandry service in Pamekasan.

**Table 1.** Dataset used in the research.

| Data Set | Number of data records | Feature | Class |
|---|---|---|---|
| Onion Data | 268 | 6 | 4 |
| Local Madura cattle | 300 | 10 | 3 |

The quality of shallots is determined by their freedom from disease, namely color conditions, freshness level, shape, weight, traces of disease, and water content. This study aims to help farmers classify the quality of onions to find more focused and optimal prospective parents. The classes built by the data are V1, ...V4, which consist of super-generative varieties of red onions as candidate seeds, super-quality varieties, medium varieties, and low varieties.

The process of selecting superior cattle seeds is based on the performance of Madura cattle body condition characteristics as a selection criterion. The quality of cattle seeds can be assessed from 10 aspects, including Healthy, normal reproductive organs, good quality, quantity of cement, Cattle Age (Months), Shoulder Height, Body length, Chest Circumference, and scrotum circumference. Superior-quality local Madura cattle: Data is taken from the types of beef cattle for breeding candidates by dividing them into 3 classes C1, C2, C3: super cattle, medium cattle, and low-quality cattle.

### 3.2 Research result

Developing an ANN requires selecting the optimal user-defined parameters, which are optimized using a large number of trial iterations. If it is known that X is a subset of q with m1 as a function of its density, and Y is also a subset of q with m2 as the density function, then a combination function m1 and m2 is m3, shown in the following equation.

**Table 2.** Setting Parameters parameters for training.

| Training cycles (Epoch) | Parametres |
|---|---|
| Neural Network | Learning rate (α)= 0.1 Momentum (m) = 0.3 Iterations= 2000 |
| ANN with Ensemble Technical | Adaboost iterations = 20, 50, 100 |
| Undersampling and Oversampling | K=100, 200 |

The model was built to overcome the class imbalance problem by combining data-level approach techniques with resampling and comparing undersampling and oversampling methods with Neural Network-based ensemble methods. Fig. 5 shows how model performance is measured by comparison.
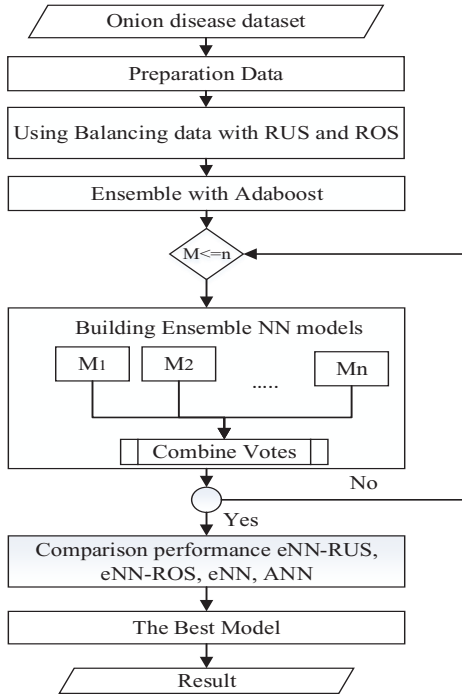


**Fig. 5.** Flow of research performance.

Several experiments were conducted using the trial parameters, namely hidden layers, by default using a learning rate value of 0.1, momentum 0.5, and a training cycle value by default until the maximum final result value. The test results are shown in Table 3 and Table 4 by displaying a comparison of several methods.

**Table 3.** Test results for several best NN variation methods on onion data.

| Approach | Acc | AUC | Kappa | F-Measure | Recall |
|---|---|---|---|---|---|
| ANN | 94.00% | 0.7582 | 0.7292 | 0.7281 | 1.0000 |
| eNN | 95.23% | 0.7871 | 0.9917 | 0.9220 | 0.8020 |
| eNN+RUS | 97.78% | 0.8739 | 0.9829 | 0.9772 | 1.0000 |
| eNN+ROS | 97.76% | 0.8629 | 0.9877 | 0.9674 | 0.9080 |

**Table 4.** Test results for several best NN variation methods on superior cattle data .

| Approach | Acc | AUC | Kappa | F-Measure | Recall |
|---|---|---|---|---|---|
| ANN | 96.00% | 0.6580 | 0.7292 | 0.7800 | 0.9800 |
| eNN | 98.05% | 0.7801 | 0.9717 | 0.8907 | 0.8849 |
| eNN+RUS | 98.00% | 0.9390 | 0.9829 | 0.9642 | 1.0000 |
| eNN+ROS | 97.69% | 0.9365 | 0.9877 | 0.9884 | 0.8608 |

The application of the sampling technique on the ensemble increased Kappa, although the accuracy decreased. Meanwhile, AUC increased when using the ensemble and the addition of sampling. On the shallot data, the performance results on eNN + RUS produced the highest performance with the highest accuracy of 97.78%, Kappa = 0.9829, and Recall = 1.0000. Meanwhile, eNN + ROS obtained the highest score on

F-Measure = 0.9776 and AUC = 0.9365, so almost the same performance was achieved with over- and under-sampling methods. Since the cow data, the results of NN + RUS produced the highest performance with the highest accuracy of 98.00%, Kappa = 0.9390, and Recall = 1.0000. Finally, adding RUS sampling produced the best performance when the primary data class was more dominant on balanced data because the data range was too high. Comparative testing of the performance of several methods showed better results with the addition of a combination of ensemble and sampling. The research shows that the ensemble's application can produce better performance than without the ensemble.
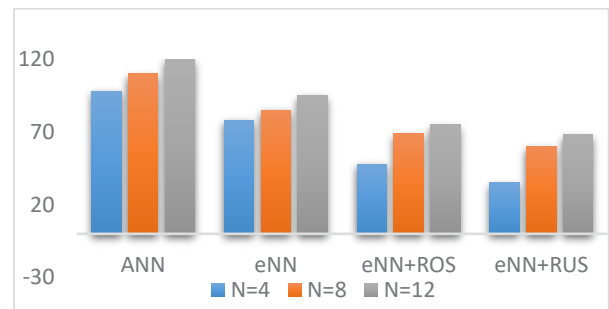


**Fig. 6.** Comparison of method execution times using the number of hidden layers using shallot data.
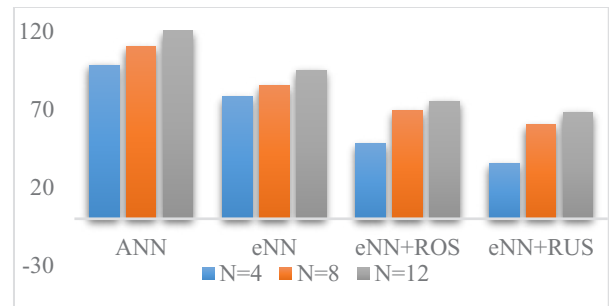


**Fig. 7.** Comparison of method execution times using the number of hidden layers using cattle data.

Fig. 6 and Fig. 7 show te more neurons used, the lowest performance will be; when n = 12, the learning process will be more rapid. Vice versa, the more Neuran is used, the more iterations will run, resulting in higher computation. The ensemble and momentum approaches show fluctuating performance depending on the best weight value, as shown in Fig. 8.
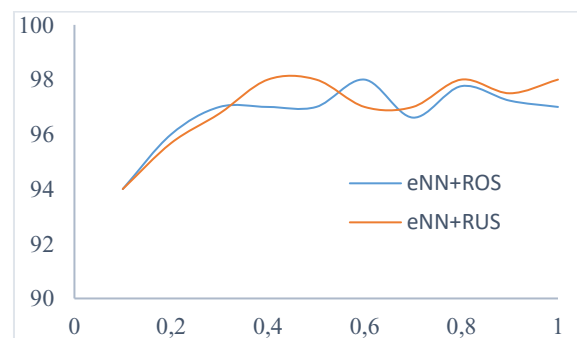


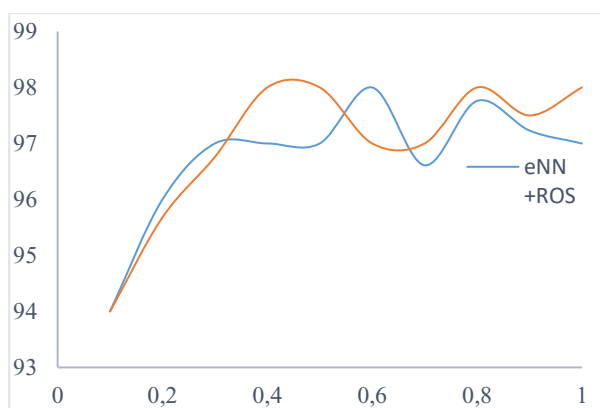**Fig. 8.** Experiment on the influence of momentum weighting on ensemble NN using shallot data.

**Fig. 8.** Experiment on the influence of momentum weighting on ensemble NN using cattle data.

ANN with momentum produces maximum accuracy when changing weights compared to without adding momentum. Increasing the momentum weight change causes significant changes in the accuracy value. The greater the given momentum value, the better the method's performance, but it is better to be less than 0.4 until 0.5 to avoid overfitting.

## 4    Conclusion

The conclusions of this research are the study's results showed that using resampling and ensemble optimization techniques on NN produced better performance in Accuracy, Kappa, and F-measure than using the basic Neural network method. Adding momentum weights to NN and using the best parameters with default iteration stops provide accuracy values close to the maximum value. Changes in weight cause quite significant changes. However, although the addition of neurons increases processing time, the processing time required will be longer by using the higher number of neurons and the additional momentum weights given. The best performance value on cattle data shows a higher accuracy value than onion data because the number of classes is smaller, and the difference in minor and significant class data is not too high, with the most considerable accuracy value of 98.05% and recall value of 1 when using the ensemble.

## References

1. J. Gong, and H. Kim, *RHSBoost: Improving classification performance in imbalance data, Computational Statistics and Data Analysis*, Elsevier B.V., **111**, 1–13, (2017), doi: 10.1016/j.csda.2017.01.005

2. Y. A. Du, *Research on the Route Pricing Optimization Model of the Car-Free Carrier Platform Based on the BP Neural Network Algorithm*, Complexity, (2021)

3. J. Choi, D. Kim, M.Ko, D. Lee, K. Wi, H. Lee, *Compressive strength prediction of ternary-blended concrete using deep neural network with tuned hyperparameters*, Journal of Building Engineering **75**, (2023)

4. M. Ahuja, D. P. Mishra, D. Mohanty, H. Agrawal, S. Roy, *Development of Empirical and Artificial Neural Network Model for the Prediction of Sorption Time to Assess the Potential of $CO_2$ Sequestration in Coal*, ACS Omega, **8,** 34, 31480-31492, (2023), https://doi.org/10.1021/acsomega.3c04412

5. S. Karmakar, G. Shrivastava, M. K. Kowar, *Impact of learning rate and momentum factor in the performance of back-propagation neural network to identify internal dynamics of chaotic motion*, Kuwait Journal of Science (KJS), **41** (2014)

6. B. Raharjo, N. Farida, P. Subekti, R. H. S. Siburian, P. D. H. Ardana, and R. Rahim, *Optimization Forecasting Using Back-Propagation Algorithm*, J. Appl. Eng. Sci., **19,** 1083–1089, (2021)

7. J. Prusa, *Using Random Undersampling to Alleviate Class Imbalance on Tweet Sentiment Data*, (2015), doi: 10.1109/IRI.2015.39

8. W. Ng. K. Dai, K. Severson, W. Huang, F. Anderson, and C. Stultz, *Generative oversampling with a contrastive variational autoencoder*, In 2019 IEEE International Conference on Data Mining (ICDM), 101–109, (2019)

9. P. Alkhairi, E. R. Batubara, R. Rosnelly, W. Wanayaumini, and H. S. Tambunan, *Effect of Gradien Descent With Momentum Backpropagation Training Function in Detecting Alphabet Letters*, Sinkron : Jurnal Penelitian Teknik Informatika, **8,** 574–583, (2023)

10. T. Akiba, S. Sano, T. Yanase, T. Ohta, M. Koyama, Optuna, *A next-generation hyperparameter optimization framework*, in Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, (2019)

11. M. G. M. Abdolrasol et al., *Artificial neural networks based optimization techniques: A review, Electron*, **10,** 21 (2021) https://doi.org/10.3390/electronics10212689

12. V. Riiman, A. Wilson, P. Pirkelbauer, *Comparing Artificial Neural Network and Cohort-Component Models for Population Forecasts*, Published in Population, Economics, Review, (2019)

13. K. M. R. Alam, N. Siddique, and H. Adeli, *A dynamic ensemble learning algorithm for neural networks*, Neural Comput. Appl., **32**, 2, 8675–8690, (2020), doi: 10.1007/s00521-019-04359-7.

14. Hairani, D. Priyanto, *A New Approach of Hybrid Sampling SMOTE and ENN to the Accuracy of*

*Machine Learning Methods on Unbalanced Diabetes Disease Data*, (IJACSA) International Journal of Advanced Computer Science and Applications, **14,** 585-590, (2023), DOI:10.14569/IJACSA.2023.0140864

15. S. S. Mullick, S. Datta, S., and S. Das, *Generative adversarial minority oversampling*. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 1695–1704, (2019)

16. Z. Huang, Y. Sang, Y. Sun, L. Jiancheng, *A neural network learning algorithm for highly imbalanced data classification*, Information Sciences, **612,** 496-513, (2022), https://doi.org/10.1016/j.ins.2022.08.074

17. C. Kasemset, K. Phuruan, T. Opassuwan, *Shallot Price Forecasting Models: Comparison among Various Techniques*, Production Engineering Archives, **29,** 4, 348-355, (2023)

18. Schlenker, Wolfram, M. J. Roberts, *Nonlinear temperature effects indicate severe damages to US crop yields under climate change*, the National Academy of sciences, **106**, 37, 15594-15598, (2009), doi: 10.1073/pnas.0906865106

19. M. Rafrin, M.Agus, P. A. Maharani, *IoT-Based Irrigation System Using Machine Learning for Precision Shallot Farming*, JURNAL RESTI (Rekayasa Sistem dan Teknologi Informasi), **8**, 2 , 216 - 222, (2024)

20. L. Sumaryanti, Nurcholis, L. Lamalewa, *Aplication of Hybrid Method for Superior cattle selection using Decision Support System*, E3S Web of Conferences, **328**, (2021).

21. Tsukazaki, Hikaru, Honjo, Masanori, Yamashita, Ken-Ichiro, Ohara, Takayoshi, Kojima, Akio, Ohsawa, Ryo, Wako, Tadayuki. *Classification and identification of bunching onion (Allium Fistulosum) varieties based on SSR markers*, Breed Sci. **60***, (2010), doi:10.1270/jsbbs.60.139.

22. Y. R. Kurniawan, D. D. Hidayat, R. Luthfiyanti, R. C. E. Andriansyah, A Indriati, *A comparative study on engineering properties of three varieties of shallots*, IOP Conf. Ser.: Earth Environ. Sci. **462**, 012025, (2020), doi:10.1088/1755-1315/462/1/012025

23. R. Benzer, *Population dynamics forecasting using artificial neural networks*, Fresenius environmental bulletin, **24**, 2, (2015)

24. A. Paul, C. Bhakat, S. Mondal, D. K. Mandal, *An observational study investigating uniformity of manual body condition scoring in dairy cows*, Indian J Dairy Sci. **73** (2020). doi:10.33785/IJDS.2020.v73i01.013

25. X. Song, E. A. M. Bokkers, S. V. Mourik, P.W.G. Groot Koerkamp, P. P. J. V. D. Tol, *Automated body condition scoring of dairy cows using 3-dimensional feature extraction from multiple body regions*, Journal of Dairy Science. **102**, 4294-4308, (2019), https://doi.org/10.3168/jds.2018-15238