SUAS Press

# Enhancing Equipment Health Prediction with Enhanced SMOTE-KNN

**ZHOU, Zhanxin** [1*] **XU, Changxin** [1] **QIAO, Yuxin** [1] **XIONG, Jize** [1] **YU, Jiqiang** [2]

[1] Northern Arizona University, USA

[2] Universidad Internacional Isabel I de Castilla, Spain

*\* ZHOU, Zhanxin is the corresponding author, E-mail: zhouzhan098@gmail.com*

**Abstract:** With the development of industrial automation, the accurate prediction of the health state of equipment becomes particularly important. This study aims to address the problem of application of small sample unbalanced data in device life prediction. A joint optimization model was constructed by combining the modified SMOTE algorithm and the modified KNN algorithm. To solve the sample imbalance problem, the modified KNN algorithm is used to improve the accuracy of classification. Through the simulation analysis of the hydraulic pump status data of Caterpillar Corporation and the vibration data of the water guide bearing of Lingjintan Hydropower Station, the proposed improved algorithm can accurately analyze the running status of the equipment and predict the future healthy development trend. Experimental results show that the joint optimization model has higher accuracy and reliability in the processing of small sample unbalanced data compared with traditional algorithms.

**Keywords:** Unbalanced Data, Equipment Life Prediction, SMOTE, Algorithm, KNN Algorithm.

## 1 Introduction

With the rapid progress of technology, China's accuracy requirement for equipment life prediction is getting higher and higher. If the key equipment fails during use, it may cause major safety accidents or economic losses. Therefore, the timely and accurate diagnosis of the health status of the device becomes an important issue. The development of machine learning and deep learning techniques provides new directions for device state prediction. Machine learning methods such as SVMs, neural networks, random forests and K nearest neighbor algorithms have been widely used in industrial production. However, these algorithms are mainly aimed at large-scale datasets and are often less effective for processing small-sample unbalanced data. Therefore, it becomes particularly important to study algorithms applicable to small-sample disequilibrium data. Although SMOTE algorithm has some advantages in handling unbalanced data as a synthetic minority oversampling technique, it still has deficiencies in near neighbor value selection, outlier processing, and unbalanced data distribution. In order to solve these problems, this study proposes an improved SMOTE algorithm (ISMOTE) and a voting KNN algorithm (VKNN) to improve the accuracy and reliability of device health state prediction by optimizing the processing of small sample unbalanced data.

## 2 Research Background

In the field of machine learning, algorithms designed for large-scale data can lead to prediction errors when handling small samples of unbalanced data, resulting in economic losses. In order to improve the prediction quality of small sample data, data enhancement is an effective method. This study uses an improved SMOTE algorithm (ISMOTE) to compensate for the deficiency of the traditional KNN algorithm in handling the imbalance and abnormal data. By the principle similar to k nearest neighbor, the ISMOTE algorithm removes scattered abnormal data and synthesqualified data while maintaining the characteristics of the data artificially. This improvement solves the problems of low quality, fuzzy boundary and abnormal distribution of the traditional SMOTE algorithm when adding new samples. Finally, the optimized data is classified by using the voting KNN algorithm (VKNN) to improve the prediction accuracy of the model.

## 3 SMOTE, The Optimization of the Algorithm

### 3.1 Application of the SMOTE Algorithm in Imbalanced Data Classification

SMOTE (Synthetic minority oversampling technology) is an improved random oversampling method

designed to solve the problem of category imbalance in data sets. Unlike simply repeatedly duplicating minority class samples, SMOTE expands the dataset by analyzing and synthesizing new minority class samples. The specific operation includes the following steps:

1) For each sample in the minority class, the Euclidean distance from the other minority class samples is calculated to determine its nearest neighbor sample.

2) Set the sampling rate according to the unbalanced ratio and determine the sampling ratio.

3) For each minority class sample, samples are randomly selected from its nearest neighbor and new samples are generated according to the specific formula:

$$x_{\text{new}} = x + rand(0,1) \cdot (\tilde{x} - x)$$

However, when device data is influenced by factors such as operating environment and device status, direct application of the SMOTE algorithm may not be suitable. The traditional SMOTE algorithm may generate new samples with noise or samples located in the boundary between the majority and minority classes, blurring the boundary of the dataset. To overcome these challenges, this paper introduces an improved SMOTE algorithm (ISMOTE), which effectively avoids the above problems and thus more accurately evaluates the health status of the device.

## 3.2 Comprehensive Optimization of Noise Filtering and Boundary

In order to solve the outliers or isolated phenomenon of the minority class samples in the majority class data group, this study proposes an improved SMOTE algorithm (ISMOTE). The algorithm evaluates each minority sample by introducing the noise proportion coefficient β, which is defined as the ratio of the number of minority samples to the number of majority samples. Using this coefficient, we can judge whether there are outliers in a specific category sample set. If the noise ratio exceeds the preset standard α, use the new sample to build the formula; otherwise, the sample point is considered as noise and removed. Based on this, this paper also proposes an optimization method for handling the problem of mixing or close distribution of normal and abnormal data in equipment data. The Euclidean distance d is calculated for each minority sample and compared with preset thresholds to optimize the distribution of minority samples. The formula is as follows:

$$d = \sqrt{\sum_{k=1}^{n} (x_{1k} - x_{2k})^2}$$

When a distance value of a minority sample point is below the predetermined threshold dmin, it is considered as an edge sample close to the majority sample sample. This method is especially suitable for the cases where the distribution of most classes and a few class samples overlaps or is abnormally sparse, and can significantly improve the classification performance and generalization ability of the data. In this paper, we use B-SMOTE algorithm and Borderline-SMOTE algorithm to classify the sample points by setting threshold and calculating distance. The optimized dataset can be efficiently applied to most modern machine learning algorithms for computation.

# 4 Application of Advanced KNN Algorithm

## 4.1 KNN Machine Learning Algorithm

The KNN algorithm, proposed by Cover and Hart in the 1960s, is an intuitive method for classification based on neighboring samples. Its core principle is that a sample roughly belongs to the same category as its nearest k neighbors. The algorithm is characterized by no prior learning of a model, but is directly classified by comparing distances between samples. Despite the simple structure of the KNN algorithm itself, scholars have been exploring how to improve its classification efficiency on particular datasets. For example, Yadav et al. studies have demonstrated the potential advantages of KNN in classification problems through mathematical analysis. Combining KNN and supersphere structures, developed the KNN-MVHM algorithm, which performs better when processing unbalanced data, Xu et al. Yin Xiaozhou et al. combined KNN with the support vector machine to form a new classifier capable of selecting the most suitable classification method under specific conditions. Li Huan et al. introduced the particle swarm optimization strategy into KNN, which improved the speed and accuracy of the nearest neighbor search. Although these improvements improve the performance of KNN algorithms, they still have limitations in dealing with overlapping and unbalanced data distribution problems, especially in applications such as device health state prediction.

## 4.2 Evolutionary KNN Classification Strategy

In this study, an evolutionary KNN (VKNN) algorithm is proposed to solve the overfitting and underfitting problems of the traditional KNN algorithm. The VKNN algorithm does not rely on the traditional majority voting principle for classification, but uses the particle swarm optimization algorithm to quickly locate the center of the sample points in the training set, and then determine a separation threshold according to the distance mean value of the center points. Finally, the sample points were "voted" and classified to achieve a more accurate classification effect. In the particle swarm optimization algorithm, the velocity and position of each particle are updated by a specific iterative formula:

$$V_{id} = \omega V_{id} + C_1 random(0,1)\left(P_{id} - X_{id}\right) + C_2 random(0,1)\left(P_{gd} - X_{id}\right)$$

$$X_{id} = X_{id} + V_{id}$$

Here, ω represents the inertia factor, and is the acceleration constant, which usually takes a value of 2. $C_1$ $C_2$ $P_{id}$ $P_{gd}$ Represents the best position found in the individual history, and represents the globally optimal position. Subsequently, the algorithm calculates and judged in the following steps:

a. Define the sample set X, containing the sample points and their class labels.

b. The fitness function f was constructed to assess the classification status of the sample points.

c. The Euclidean distance of the sample point to the center of its category was calculated to obtain the distance set D.

d. The distance mean was determined and used as a threshold to divide the sample points.

e. According to the mean distance, decide whether the sample points belong to the isolated sample set, and judge the sample points in the category. $D_{new}$ $D_{new}$

In this way, VKNN algorithm solves the problem of inconsistent classification results caused by traditional KNN with different k values, and improves the computational efficiency. In practice, this improvement means that the computing time required can be significantly reduced while maintaining high accuracy, which is important for improving productivity.

# 5 Implementation Strategy of the Integrated ISMOTE-VKNN Algorithm

In this study, the ISMOTE-VKNN algorithm was designed to optimize the classification of device status data. As shown in Figure 1, first, the data are cleaned through the preprocessing phase and divided into the training set and the test set in a ratio of approximately 2:1. Then, the noise coefficient β for each sample point is calculated, based on the set noise threshold, α, to pick out the samples suitable for model training. Next, by setting a distance threshold of $d_{min}$, To out suitable sample points within d. These sample points were subjected to proximity values analysis to identify an suitable sample generation strategy. In addition, the algorithm uses particle swarm optimization (PSO) technology to locate the center point of each category of sample, and calculate the distance d based on this center point, and then separate the sample points and generate a new sample set $D_{new}$. Ultimately, this new sample set was classified using the voting mechanism to produce the final classification results.
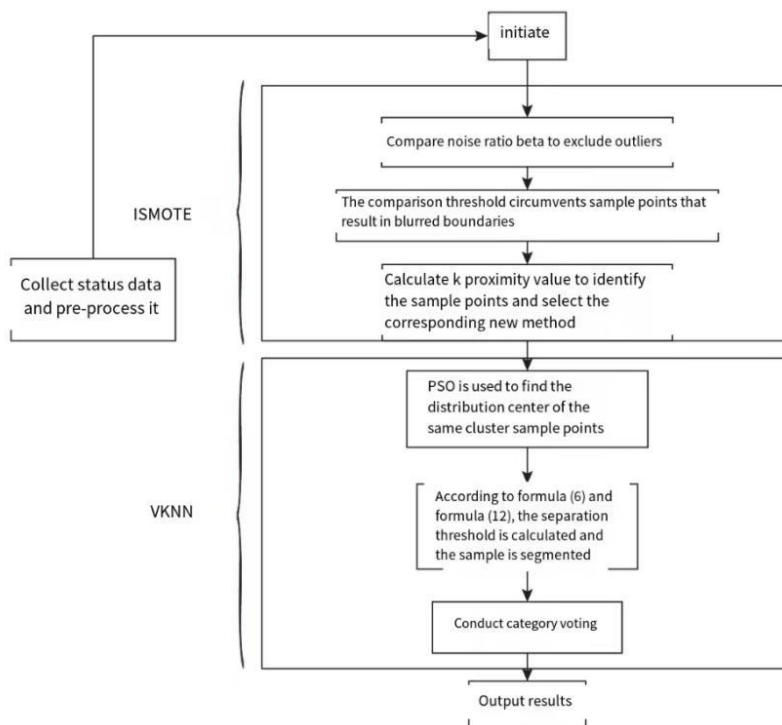


**Figure 1 ISMOTE-Execution procedure of the VKNN algorithm**

# 6 Experimental Validation of the ISMOTE-VKNN Algorithm on Device Running Data

## 6.1 Selection of the Data

To verify the effectiveness of ISMOTE-VKNN algorithm, this study utilized the hydraulic pump data provided by Caterpillar and the operation data of water guide bearing of unit 8 of Lingjintan Hydropower Station. The vibration data for these devices contain important information about potential failures, often reflected in changes in bearing vibration. In the hydraulic pump case, a continuous 1-minute data acquisition was performed at 10-minute intervals, followed by a feature extraction to apply to the model used in this study. For the water guide bearing data, the vibrations at different working loads are recorded for the analysis of their lifetime. During the experiment, the first two thirds of the hydraulic pump data were used to train the model, and the rest was used to test the predictive performance of the model. The water guide bearing data was processed and allocated in the same proportion as the hydraulic pump data to ensure the accuracy of the experimental results. All the experiments were performed in the Anaconda 3.0 environment.

## 6.2 2-dimensional Transformation and Preprocessing Effect of Vibration Data

In this study, to simplify the analysis of multi-dimensional vibrational data, it was converted into two-dimensional form and representative datassets CH 1-1 and CH 1-9 were selected for simulation experiments. As shown in Figure 2, the processed data demonstrate the deviation of a few class sample points in the original dataset, namely isolated outlier points, which are shown in black in the data distribution map, while the normal running state data points are shown in orange. The processing of such data directly using conventional KNN algorithms is prone to prediction bias.
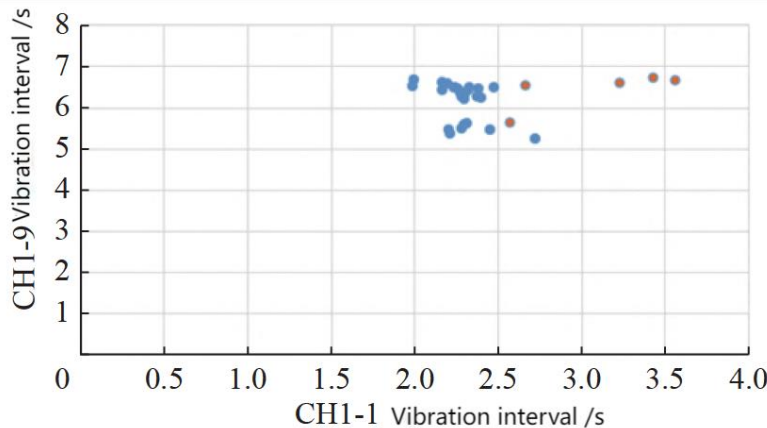


**Figure 2. Distribution of the vibration data**

## 6.3 Application of Optimization Processing in Outlier Identification

In the experimental phase, the best k value of 4 was determined by a Bayesian posterior probability analysis. Next, the noise ratio in the minority class samples was calculated by setting a specific threshold, used to identify and exclude outliers in the data. The selected noise ratio threshold $\alpha$ was set to 0.1 to facilitate clear discrimination of outliers, as shown in Figure 3, data points with a $\beta$ value 0 were treated as unqualified and removed from the dataset. Although outliers were excluded, there may still be some data points too close between majority and minority samples. To prevent these points from causing blurred classification boundaries, a more stringent distance threshold d was set$_{min}$Was 0.5, and the sample points were screened according to this threshold. In the ISMOTE-KNN model, the data were further processed by setting m=0.5 and k=4, and the eligible samples were randomly oversampled. After oversampling, a new set of sample points was obtained. To ensure that the distribution characteristics of the raw data remained constant, the lifetime curves of the raw and new data points were fitted. The results show that the lifetime fitting curves of the original data are almost unaffected even after adding new data points, demonstrating the effectiveness and feasibility of the ISMOTE algorithm for data optimization.
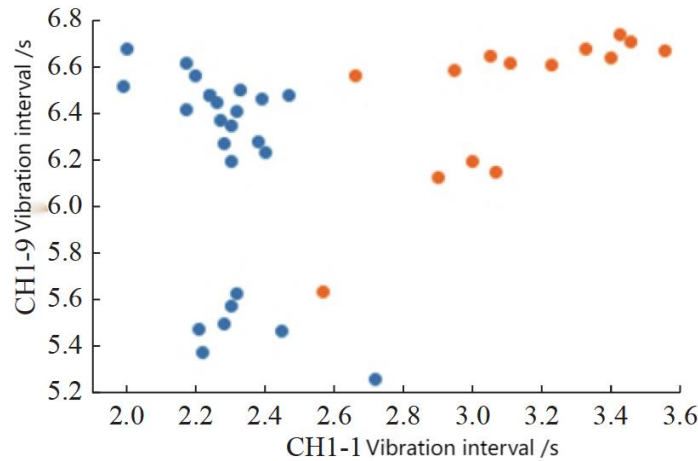
**Figure 3 ISMOTE-KNN, the distribution of the vibration data for the model**

## 6.4 Study of Data Processing and Classification Efficiency of VKNN

In this study, the evolutionary KNN (VKNN) algorithm was used to classify the vibration data optimized by the ISMOTE algorithm. With the help of particle swarm optimization (PSO) algorithm, the parameters of VKNN algorithm are set as follows: the inertia factor ω is gradually reduced from the initial value 0.9 to the end value 0.4, and the two acceleration constants C1 and C2 are set to 1.5. In this configuration, when the maximum number of iterations reaches 100, the PSO successfully determines the center

point of the same cluster sample at (3.12,6.45) and (2.29,6.16). The calculated distance threshold dmin was 0.373 and 0.411, respectively, providing evidence for the classification of the data. Subsequently, the final classification result was determined through the "voting" mechanism. As shown in the experimental results, after applying the VKNN algorithm to the 12 test data, the accuracy of device health state prediction increased significantly, and the accuracy reached 94.9% in the training set and 100% in the test set. Compared with the direct application of KNN algorithm, ISMOTE-VKNN algorithm demonstrated significant advantages of speed and accuracy in processing small sample data.

**Table 1 Comparison of accuracy of different algorithms on the hydraulic pump dataset**

| The algorithm category | Algorithm name | Training set accuracy (%) | Test set accuracy (%) |
|---|---|---|---|
| traditional algorithm | tradition KNN | 93.3 | 91.7 |
| traditional algorithm | Conventional nonlinear SVM | 90.0 | 66.7 |
| Joint algorithm | ISMOTE_SVM | 92.3 | 66.7 |
| Joint algorithm | ISMOTE_VKNN | 94.9 | 100 |

To further evaluate the performance of the ISMOTE-VKNN algorithm, this study was compared with the nonlinear support vector machine (SVM) algorithm. The results showed that the nonlinear SVM algorithm increased from 90.0% to 92.3% on the training set, while the accuracy on the test set remained unchanged. This shows that although the ISMOTE algorithm does not show significant advantages on the test set, the ISMOTE algorithm is more effective on the training set with a larger amount of data.

**Table 2 Comparison of the classification accuracy of the SVM and the improved SVM**

| Data category | The SVM raw data accuracy (%) | Improved SVM ISMOTE _ SVM accuracy rate (%) |
|---|---|---|
| training set | 90.0 | 92.3 |
| test set | 66.7 | 66.7 |

When analyzing the small sample data of the hydraulic pump, although the error rate of the training set decreased, the error rate of the test set was still high, similar to the results obtained directly using the KNN algorithm. This observation further highlights the particular advantage of the

ISMOTE-VKNN algorithm in handling small samples of data. Moreover, when the same method is applied to water guide bearing vibration data, the results show that the ISMOTE-VKNN algorithm has higher classification accuracy in real scenarios compared to traditional machine

**Journal of Industrial Engineering and Applied Science**
**ISSN 3005-6071 (Print) | ISSN 3005-608X (Online) | Vol. 2, No. 2, 2024**

SUAS Press

learning algorithms, even when a few classes have fewer samples.

# 7 Conclusion

In this study, we propose a joint optimization model combining improved SMOTE algorithm with the application of small sample unbalanced data in improved KNN algorithm for device life prediction. The processing and analysis of experimental data demonstrate the effectiveness of the ISMOTE-VKNN algorithm in improving classification accuracy and processing small sample data. The experimental results show that the proposed algorithm can predict the health status of the device more accurately compared with the traditional KNN and SVM algorithms, especially in small sample cases. Future studies could further explore the application of the algorithm in other fields and how to optimize the algorithm parameters to accommodate a wider range of data features to improve its universality and reliability in practical industrial applications.

# Institutional Review Board Statement

Not applicable.

# Informed Consent Statement

Not applicable.

# Data Availability Statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

# Conflict of Interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

# Publisher's Note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

# Author Contributions

Not applicable.

# About the Authors

**ZHOU, Zhanxin**

Affiliation: Northern Arizona University.

**XU, Changxin**

Affiliation: Northern Arizona University.

**QIAO, Yuxin**

Affiliation: Northern Arizona University.

**XIONG, Jize**

Affiliation: Northern Arizona University.

**YU, Jiqiang**

Affiliation: Universided Internacional Isabel I of Castile.

# References

[1] Liu, T., Xu, C., Qiao, Y., Jiang, C., & Chen, W. (2024). News Recommendation with Attention Mechanism. Journal of Industrial Engineering and Applied Science, 2(1), 21-26.

[2] Liu, T., Xu, C., Qiao, Y., Jiang, C., & Yu, J. (2024). Particle Filter SLAM for Vehicle Localization. Journal of Industrial Engineering and Applied Science, 2(1), 27-31.

[3] Xu, C., Qiao, Y., Zhou, Z., Ni, F., & Xiong, J. (2024). Accelerating Semi-Asynchronous Federated Learning. arXiv preprint arXiv:2402.10991.

[4] Zhou, Z. (2024, February). ADVANCES IN ARTIFICIAL INTELLIGENCE-DRIVEN COMPUTER VISION: COMPARISON AND ANALYSIS OF SEVERAL VISUALIZATION TOOLS. In The 8th International scientific and practical conference "Priority areas of research in the scientific activity of teachers"(February 27–March 01, 2024) Zagreb, Croatia. International Science Group. 2024. 298 p. (p. 224).

[5] Xu, C., Yu, J., Chen, W., & Xiong, J. (2024, January). DEEP LEARNING IN PHOTOVOLTAIC POWER GENERATION FORECASTING: CNN-LSTM HYBRID NEURAL NETWORK EXPLORATION AND RESEARCH. In The 3rd International scientific and practical conference "Technologies in education in schools and universities"(January 23-26, 2024) Athens,

**Journal of Industrial Engineering and Applied Science**
**ISSN 3005-6071 (Print) | ISSN 3005-608X (Online) | Vol. 2, No. 2, 2024**

SUAS
Press

Greece. International Science Group. 2024. 363 p. (p. 295).

[6] Su, J., Jiang, C., Jin, X., Qiao, Y., Xiao, T., Ma, H., ... & Lin, J. (2024). Large Language Models for Forecasting and Anomaly Detection: A Systematic Literature Review. arXiv preprint arXiv:2402.10350.

[7] Qiao, Y., Jin, J., Ni, F., Yu, J., & Chen, W. (2023). Application of machine learning in financial risk early warning and regional prevention and control: A systematic analysis based on shap. WORLD TRENDS, REALITIES AND ACCOMPANYING PROBLEMS OF DEVELOPMENT, 331.

[8] Qiao, Y., Ni, F., Xia, T., Chen, W., & Xiong, J. (2024, January). AUTOMATIC RECOGNITION OF STATIC PHENOMENA IN RETOUCHED IMAGES: A NOVEL APPROACH. In The 1st International scientific and practical conference "Advanced technologies for the implementation of new ideas"(January 09-12, 2024) Brussels, Belgium. International Science Group. 2024. 349 p. (p. 287).

[9] Ni, F., Zang, H., & Qiao, Y. (2024, January). SMARTFIX: LEVERAGING MACHINE LEARNING FOR PROACTIVE EQUIPMENT MAINTENANCE IN INDUSTRY 4.0. In The 2nd International scientific and practical conference "Innovations in education: prospects and challenges of today"(January 16-19, 2024) Sofia, Bulgaria. International Science Group. 2024. 389 p. (p. 313).

[10] Su, J., Nair, S., & Popokh, L. (2022, November). Optimal Resource Allocation in SDN/NFV-Enabled Networks via Deep Reinforcement Learning. In 2022 IEEE Ninth International Conference on Communications and Networking (ComNet) (pp. 1-7). IEEE.

[11] Zhao, Z., Zhang, N., Xiong, J., Feng, M., Jiang, C., & Wang, X. (2024). Enhancing E-commerce Recommendations: Unveiling Insights from Customer Reviews with BERTFusionDNN. Journal of Theory and Practice of Engineering Science, 4(02), 38-44.

[12] Xiong, J., Feng, M., Wang, X., Jiang, C., Zhang, N., & Zhao, Z. (2024). Decoding sentiments: Enhancing covid-19 tweet analysis through bert-rcnn fusion. Journal of Theory and Practice of Engineering Science, 4(01), 86-93.

[13] Bao, W., Che, H., & Zhang, J. (2020, December). Will_Go at SemEval-2020 Task 3: An accurate model for predicting the (graded) effect of context in word similarity based on BERT. In Proceedings of the Fourteenth Workshop on Semantic Evaluation (pp. 301-306).

[14] Su, J., Nair, S., & Popokh, L. (2023, February). EdgeGym: A Reinforcement Learning Environment for Constraint-Aware NFV Resource Allocation. In 2023

IEEE 2nd International Conference on AI in Cybersecurity (ICAIC) (pp. 1-7). IEEE.

[15] Su, J., Nair, S., & Popokh, L. (2022, November). Optimal Resource Allocation in SDN/NFV-Enabled Networks via Deep Reinforcement Learning. In 2022 IEEE Ninth International Conference on Communications and Networking (ComNet) (pp. 1-7). IEEE.

[16] Huang, X., Zhang, Z., Guo, F., Wang, X., Chi, K., & Wu, K. (2024). Research on Older Adults' Interaction with E-Health Interface Based on Explainable Artificial Intelligence. arXiv preprint arXiv:2402.07915.

[17] Zhang, X., Guo, F., Chen, T., Pan, L., Beliakov, G., & Wu, J. (2023). A Brief Survey of Machine Learning and Deep Learning Techniques for E-Commerce Research. Journal of Theoretical and Applied Electronic Commerce Research, 18(4), 2188-2216.

[18] Li, X., Ye, W., Zou, M., Guo, F., & Huang, Z. (2023, October). Analysis and research of retail customer consumption behavior based on support vector machine. In 2023 IEEE 3rd International Conference on Data Science and Computer Application (ICDSCA) (pp. 713-718). IEEE.

[19] Zhang, Z., Li, P., Hammadi, A. Y. A., Guo, F., Damiani, E., & Yeun, C. Y. (2023). Reputation-Based Federated Learning Defense to Mitigate Threats in EEG Signal Classification. arXiv preprint arXiv:2401.01896.

[20] Peng, Y., Bian, J., & Xu, J. (2024). FedMM: Federated Multi-Modal Learning with Modality Heterogeneity in Computational Pathology. arXiv preprint arXiv:2402.15858.

[21] Wang, L., Bian, J., & Xu, J. (2023). Federated Learning with Instance-Dependent Noisy Labels. arXiv preprint arXiv:2312.10324.

[22] Bian, J., Ren, S., & Xu, J. (2023). CAFE: Carbon-Aware Federated Learning in Geographically Distributed Data Centers. arXiv preprint arXiv:2311.03615.

[23] Bian, J., & Xu, J. (2023, October). Client Clustering for Energy-Efficient Clustered Federated Learning in Wireless Networks. In Adjunct Proceedings of the 2023 ACM International Joint Conference on Pervasive and Ubiquitous Computing & the 2023 ACM International Symposium on Wearable Computing (pp. 718-723).

[24] Bian, J., Shen, C., & Xu, J. (2023). Joint Client Assignment and UAV Route Planning for Indirect-Communication Federated Learning. arXiv preprint arXiv:2304.10744.

[25] Xu, J., & Sen, S. (2023). Ensemble variance reduction methods for stochastic mixed-integer programming and their application to the stochastic facility location problem. INFORMS Journal on

Computing.

[26]    Xu, J., & Sen, S. (2023). Compromise policy for multi-stage stochastic linear programming: Variance and bias reduction. Computers & Operations Research, 153, 106132.

[27]    Liu, S., Wu, K., Jiang, C., Huang, B., & Ma, D. (2023). Financial time-series forecasting: Towards synergizing performance and interpretability within a hybrid machine learning approach. arXiv preprint arXiv:2401.00534.

[28]    Wei, K., Zang, H., Pan, Y., Wang, G., & Shen, Z. (2024). Strategic application of ai intelligent algorithm in network threat detection and defense. Journal of Theory and Practice of Engineering Science, 4(01), 49-57.

[29]    Zang, H., & Dong, X. (2024). Optimizing Soil Health Management in Smart Agriculture: Deep Learning Algorithms for Nutrient Analysis and Fertilizer Recommendation with Precision Agriculture Systems. Journal of Industrial Engineering and Applied Science, 2(1), 1-7.

[30]    Zang, H. (2024). Precision calibration of industrial 3d scanners: An ai-enhanced approach for improved measurement accuracy. Global Academic Frontiers, 2(1), 27-37.