

Interpretability of Unsupervised Anomaly Detection Model: One-Class Support Vector Machine with Rule Extraction

Ji Hun Park, Sang Won Oh, Man Gyun Na*

Department of Nuclear Engr., Chosun Univ., 10 Chosundae 1-gil, Dong-gu, Gwangju, Republic of Korea, 61452

*Corresponding author: magna@chosun.ac.kr

***Keywords:** interpretability, unsupervised model, anomaly detection, one-class SVM, rule extraction

1. Introduction

Nuclear Power Plants (NPPs) are safety-critical facilities, and minimizing human errors during operation is essential for maintaining safety. Accordingly, many studies are being conducted to reduce human errors. Among them, Artificial Intelligence (AI)-based decision support systems are being researched to minimize human errors by operators. However, the black-box characteristics of AI need to be addressed before these studies can be applied in the field.

Recently, eXplainable AI (XAI) methods have been developed to address the black-box characteristics of AI systems. The XAI methods can provide users with an explanation of why the AI's output values were derived. These XAI methods are being applied to AI-based decision support systems to reduce human error in NPPs. AI-based decision support systems are categorized into supervised learning-based classification and regression problems, as well as unsupervised learning-based anomaly detection. The application of XAI to supervised learning-based models is an active area of research. On the other hand, the application of XAI to unsupervised learning-based AI models is still limited. Therefore, this study utilizes the One-Class Support Vector Machine (OCSVM) with rule extraction method. It combines OCSVM, an unsupervised learning-based AI method, and rule extraction method, an XAI method. As part of the primary research to evaluate the applicability of the OCSVM with rule extraction, two case studies are conducted as follows:

- (1) Anomaly detection using accelerated aging data of Insulated Gate Bipolar Transistor (IGBT)
- (2) Anomaly detection using NPP simulation data

The IGBT data are sourced from open-source data provided by the National Aeronautics and Space Administration (NASA), and NPP simulation data are collected using a Compact Nuclear Simulator (CNS). An anomaly detection model is developed and evaluated in the case study using the OCSVM method. Additionally, rule extraction is performed based on the pre-trained anomaly detection model. The extracted rules can identify the logical structure of OCSVM (e.g., the boundary of the normal region can be identified as the rules). Consequently, it is expected to be a useful tool to enhance the interpretability of OCSVM-based anomaly detection models.

2. Methods

This section introduces OCSVM and rule extraction and describes the method that combines them, known as OCSVM with rule extraction.

2.1 OCSVM

The OCSVM method [1] is derived from traditional SVM method. The SVM method is a supervised learning algorithm primarily utilized for classification and regression problems. As a supervised learning algorithm, it necessitates a labeled dataset. Additionally, it employs a kernel function to transform the data into a high-dimensional feature space. This enables it to solve non-linear problems and facilitates the recognition of relationships between the data. The primary goal of the SVM method is to establish a decision boundary between datasets. Fig. 1 illustrates the hyperplane and the line with support vectors separating two datasets (i.e., '+' and '-'). Here, the greater the distance between the hyperplane and the line with support vectors, the higher the confidence level, indicating that training progresses toward enlarging the margin.

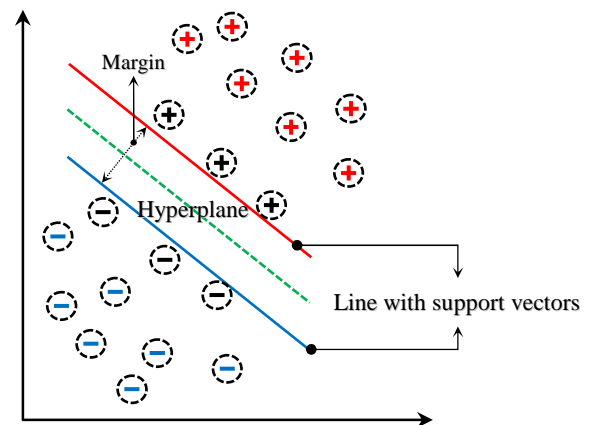


Fig. 1. SVM configuration; hyperplane, line with support vectors, and margin.

The OCSVM method is trained similarly to the SVM method. However, the OCSVM method differs because it is an unsupervised learning algorithm. This means that unlabeled data are used; thus, the margin between the data is not utilized, as depicted in Fig. 1. In the OCSVM

method, the margin is the distance between the origin and the hyperplane, as depicted in Fig. 2; similarly to the SVM method, training aims to maximize the margin.

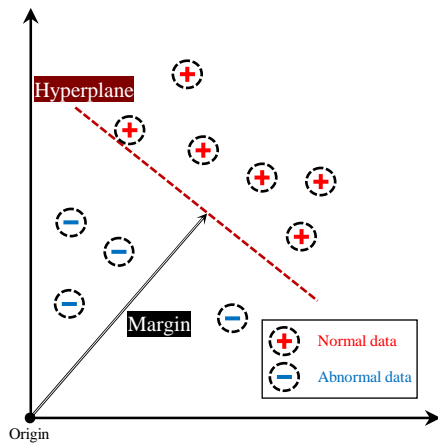


Fig. 2. OCSVM configuration; hyperplane, and margin.

The primary hyper-parameters of the OCSVM method are ν and γ . ν represents the upper bound on the learning error rate and the lower bound on the support vector rate (which primarily determines the anomaly detection boundary). γ is the coefficient for the radial basis function kernel, influencing the curvature of the decision boundary. The complexity of the algorithm increases as both hyper-parameter values increase, so it is crucial to find the parameter values that optimized performance.

2.2 OCSVM with rule extraction

The rule extraction method [2] employs a clustering algorithm to generate a hypercube. The hypercube represents boundary regions in a multi-dimensional feature space. The implementation of this hypercube involves the following steps:

- (1) Obtaining clustering results:
 - Clustering algorithms such as k-means and k-prototypes are employed to group the data into clusters, utilizing the outcomes of the OCSVM-based anomaly detection model.
- (2) Finding the boundary points of each cluster:
 - To determine the boundaries of each cluster, certain data points belonging to that cluster are selected.
 - The minimum and maximum values of the selected data points define the boundaries, and this process is repeated for each variable.
- (3) Constructing a hypercube in a multi-dimensional feature space:
 - Combining the ranges for each variable to create a hypercube.
- (4) Verifying the hypercube boundaries:

- Verify if outliers exist within the hypercube and repeat the hypercube creation process until no outliers are present.

The boundary conditions of these hypercubes serve as rules. In other words, the boundary conditions of the OCSVM-based anomaly detection model can be articulated as rules.

3. Data Preparation

This section describes the datasets used in the case study and discusses data standardization.

3.1 IGBT accelerated aging data

The data on accelerated aging of IGBTs are obtained from open-source data [3]. The IGBT data were subjected to accelerated aging by applying temperature and voltage conditions higher than the design characteristics. As the IGBT ages, it will fail, with the failure symptom being latch-up. Latch-up causes a low output voltage relative to the supply voltage, as shown in Fig. 3. The difference between the supply and output voltage is defined as the degradation characteristic. The failure criterion is the point where the output voltage drops sharply. The IGBT datasets include information for 4 devices, experimental temperature and voltage, and degradation characteristics. For anomaly detection model development, we only use operation time, temperature, and voltage. The other variables, including degradation characteristics, are not available in the field and are therefore excluded from the AI model development. Additionally, we used normal condition data for 3 devices to train the anomaly detection model. The normal condition data for the remaining 1 device and all anomaly data were used for validation.

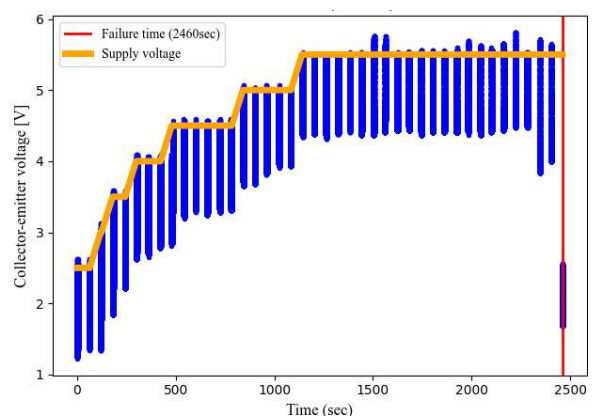


Fig. 3. Degradation characteristic of IGBT data; note the low output voltage (blue line) relative to the supply voltage (orange line).

3.2 NPP simulation data

The NPP simulation data were collected using the CNS. The collected data were divided into normal operating data for training the anomaly detection model and abnormal operating data for validation. Loss of coolant accident and steam generator tube rupture scenarios are utilized for the abnormal operating data. Similar to the IGBT dataset, only the normal operating data are used for training. The validation and test data utilize a subset of the normal operating data and all of the abnormal operating data. 137 variables were selected as input for the AI based on the symptom requirements of each scenario.

3.3 Data standardization

Data standardization is a scaling method that makes the mean 0 and the variance 1 for each variable; that is, the values are transformed to have a Gaussian normal distribution. This is expressed as in Eq. (1).

$$x' = \frac{x - \text{mean}(x)}{\text{std}(x)} \quad (1)$$

This not only prevents learning from being dependent on the magnitude of the variable values, but also contributes to faster learning.

4. Case Study

This section discusses a series of case studies that utilize the OCSVM method to develop anomaly detection models and then extract rules. The used data are IGBT data and NPP simulation data.

4.1 Development of anomaly detection model and rule extraction using IGBT data

An anomaly detection model is developed and optimized using the IGBT data. Then, based on the pre-trained anomaly detection model, rules are extracted using a rule extraction method. The training data are normal data for 3 devices, as described earlier. The validation and test data consist of normal and anomaly for the remaining 1 device. A grid search algorithm is utilized, and the ν and γ values are optimized. The optimized model is selected based on three metrics calculated using the values (i.e., TP, FP, FN, TN) from Table I.

Table I: Confusion Matrix for Anomaly Detection Model Evaluation

Predicted	Expected	
	Normal	Abnormal
Normal	TP	FP
Abnormal	FN	TN

The optimization criteria are based on the following evaluation metrics:

- (1) Metric 1: percentage of samples predicted to be normal that are actually normal (for training data; normal data only) (refer to Eq. (2))

$$\text{Metric1} = \frac{TP}{TP+FP} \text{ for training data} \quad (2)$$

- (2) Metric 2: percentage of samples predicted to be normal that are actually normal (for validation data; normal data only) (refer to Eq. (3))

$$\text{Metric2} = \frac{TP}{TP+FP} \text{ for validation data} \quad (3)$$

- (3) Metric 3: percentage of samples predicted to be abnormal that are actually abnormal (for validation data; abnormal data only) (refer to Eq. (4))

$$\text{Metric3} = \frac{TN}{TN+FN} \text{ for validation data} \quad (4)$$

The highest performance is observed for the hyper-parameter condition with $\nu=0.1$ and $\gamma=0.4$; the performance percentages for each metric are 90.57% (metric 1), 96.9% (metric 2), and 100% (metric 3). The results obtained through rule extraction based on the optimized model are shown below, with 9 rules were extracted. The rules are divided into the ranges of the input variables, such as time (operation time), temperature, and voltage.

- (1) ($62 \leq \text{time} \leq 784$) and ($100 \leq \text{temperature} \leq 240$) and ($2.5 \leq \text{voltage} \leq 4.5$)
- (2) ($902 \leq \text{time} \leq 1504$) and ($265 \leq \text{temperature} \leq 280$) and ($5 \leq \text{voltage} \leq 5.5$)
- (3) ($600 \leq \text{time} \leq 1085$) and ($240 \leq \text{temperature} \leq 265$) and ($\text{voltage} = 5$)
- (4) ($2280 \leq \text{time} \leq 2281$) and ($\text{temperature} = 100$) and ($\text{voltage} = 5.5$)
- (5) ($1440 \leq \text{time} \leq 1805$) and ($265 \leq \text{temperature} \leq 280$) and ($5 \leq \text{voltage} \leq 6$)
- (6) ($1861 \leq \text{time} \leq 1924$) and ($\text{temperature} = 280$) and ($\text{voltage} = 6$)
- (7) ($1800 \leq \text{time} \leq 1985$) and ($265 \leq \text{temperature} \leq 280$) and ($\text{voltage} = 5.5$)
- (8) ($2041 \leq \text{time} \leq 2223$) and ($\text{temperature} = 265$) and ($\text{voltage} = 5.5$)
- (9) ($2042 \leq \text{time} \leq 2105$) and ($\text{temperature} = 280$) and ($\text{voltage} = 5.5$)

4.2 Development of anomaly detection model and rule extraction using CNS data

An anomaly detection model is developed and optimized using CNS data. Similarly to the IGBT data, rules are extracted based on a pre-trained model. The

data were split into normal and abnormal data, as described in subsection 3.2, to serve as training, validation, and test data. The performance of the anomaly detection model using CNS data is as follows: 95.12% (metric 1), 97.35% (metric 2), and 99.6% (metric 3). The hyper-parameters are $\nu=0.1$ and $\gamma=0.6$. The extracted rules based on the optimized model are shown in Fig. 4. All 137 variables are used as building blocks for the rules. As a result, 86 rules were extracted.

```
*****
Combination of categorical variables N°1
----- Subgroup 1 -----
Rule N° 1: IF BFPV145 = 0.646632731277466 AND ZCNDTR = 0.9477934319761384 AND CIOADMPC = 99.29454803466795 AND UCHGUT = 280.2807861238125 AND BLV614 = 0.0103086691094448 AND QPRZH = 0.94039002418518 AND ZREAC = 6.4624671936035165 AND ZNST63 = 55.14003735661098 AND ZNST76 = 66.9437879023363 AND WFLWNS = 156.3467869202344 AND ZVCT = 54.33377016511 AND BHTV = 0.9730818867683411 AND CAXOFF = 0.5235962629447931 AND UPRZ = 340.31561279296875 AND CRETV = 1.720226947546927407 AND ZNST85 = 534.489844373 AND ZNST76 = 2.53299474761865 AND ZNST87 = 528.1591798875 AND WSRFRAY = 0.0 AND BHTV = 1.0 AND BHRVCV = 0.9905047414016782 AND ZNST76 = 0.0 AND ZNST103 = 0.17140188936633 AND ZNST76 = 534.4648003322 AND CIOADMPC = 2663.744106625 AND UREDUT = 152.5839221411094 AND ZCND = 0.778953430368042 AND ZNST108 = 39.87906845125734 AND ZNST88 = 1.00697535148813 AND ZNST1 = 99.277610778806 AND BPRCV = 0.0 AND BFPV122 = 0.5474188923835754 AND WBOAC = 0.0 AND ZVCT = 0.9430394717911 AND UPRZ = 34.670367590483 AND UALLEG = 909.091247358975 AND BTV18 = 0.0 AND ZNST74 = 0.0 AND ZNST81 = 98.3774642044336 AND ZNST67 = 1.006589651107879 AND ULPHOT = 0.778953430368042 AND ZNST108 = 155.545166015625 AND ZNST72 = 88.75740202641016 AND ZNST78 = 0.0026438120978 AND ZNST69 = 1.0349740791512079 AND UPRZ = 147.944380078125 AND QOVR = 1.0873954910783 AND UPRZ = 336.59632733906 AND ZNST101 = 99.623930333964 AND KBCD015 = 15.0 AND ZNST79 = 100.0849699696795 AND KBCD07 = 20.0 AND BFPV488 = 0.924419710995938 AND WFLWNS = 534.562072739061 AND UALLEG = 100.5338652359061 AND BFPV479 = 0.01830245584204 AND QPRZH = 40.00765991109575 AND ZNST100 = 236.32865907617 AND ZNST74 = 65.314607343943 AND ZNST79 = 66.7823394217811 AND UPRZ = 487.9818463253 AND ZVCT = 53.15314869869067 AND ZNST80 = 100.0903470900906 AND UALLEG = 326.596822753906 AND ZNST101 = 88.7230289169922 AND UCHGNO = 4.399923709106445 AND ZNST121 = 51.410105 AND BPRV10 = 0.0 AND UALLEG = 390.4804887 AND ZNST11 = 0.2412179782651 AND ZNST10 = 8.919460575836 AND UALLEG = 399.3861030274575 AND BFPV48 = 0.87184343593594 AND WFLWNS = 31.48281683793954 AND QPRZLD = 0.9297760263212 AND ZNST70 = 88.7741897830078 AND ZNST73 = 65.31695556462625 AND ZNST124 = 99.6872682324216 AND ZNST76 = 0.0022894207 AND UPRZ = 34.670367590483 AND ZNST102 = 9.999999473247407 AND BLV145 = 1.0 AND BFPV49 = 0.52118328224779 AND UALLEG2 = 292.1361093894375 AND BFPV489 = 0.43492591810173 AND ZNST77 = 50.03113174438477 AND ZNST75 = 65.308465860316 AND UALLEG = 309.3670049453125 AND EBOAC = 99998.78125 AND UALLEG3 = 326.596822753906 AND UALLEG3 = 292.131459969975 AND ZAPWTK = 8.71857929297363 AND UALLEG = 909.091247358975 AND BTV18 = 0.0 AND ZNST74 = 0.0 AND ZNST81 = 0.294590471504883 AND FROEN = 60.0 AND PVAC = 721.2213134765625 AND KBCD022 = 99.0 AND BHTV2 = 1.0 AND PVCT = 1.733920469789425 AND KBCD011 = 160.0 AND ZNST36 = 29.253202078125 AND UALLEG3 = 309.3218078613281 AND BHTV6 = 1.0 AND ZNST3 = 0.0 AND KBCD019 = 1800.0 AND KBCD06 = 228.0 AND ZNST6 = 0.231696079964061 AND KBCD06 = 402.0 AND BHTV1 = 0.0 AND FSRMDPM = 1.68153817189780544 AND HCOIC = 0.0 AND BHTV22 = 0.0 AND BHTV30 = 0.0 AND BHTV13 = 0.0 AND KBCD08 = 228.0 AND ZNST22 = 2.0791811948342 AND BHTV = 1.0 AND KLAMPO15 = 1.0 AND BLV48 = 0.0 AND KBCD05 = 228.0 AND KBCD020 = 998.0 AND KBCD010 = 228.0 AND BTV143 = 1.0 AND KLAMPO48 = 0.0 AND KBCD04 = 228.0 AND KLAMPO21 = 1.0 AND KLAMPO24 = 1.0 AND BHTV1 = 0.0 AND BFPV145 = 0.64619047176671 AND ZCNDTR = 9.474410057067873 AND CIOADMPC = 99.29454803466795 AND UCHGUT = 272.3257141113207 AND BLV614 = 0.0 AND QPRZH = 0.870232820877052 AND ZREAC = 6.4624671936035165 AND ZNST63 = 55.14003735661098 AND ZVCT = 52.124298734839 AND ZNST85 = 156.30016479492188 AND ZVCT = 73.89144134514884 AND BHTV = 0.972910940641222 AND CAXOFF = 0.524028420448302 AND UPRZ = 340.17596435548875 AND CRETV = 5.76533878415930468 AND ZNST85 = 534.4400376171875 AND ZNST56 = 2.389007806777954 AND ZNST87 = 528.1431884765625 AND WSRFRAY = 0.0 AND BHTV1 = 1.0 AND BHRVCV = 0.995976388444371 AND ZNST66 = 0.0 AND ZNST103 = 87.148597712816 AND ZNST86 = 534.419616692189 AND CIOADMPC = 2663.73101021875 AND UREDUT = 142.681684011877 AND ZCND = 0.778601646423598 AND ZNST108 = 39.87906845125734 AND WBOAC = 1.006656271912929 AND ZNST1 = 99.2520338330078 AND BPRV = 0.0 AND BHTV12 = 0.949208199977875 AND WBOAC = 0.542174147603896 AND WFLWNS = 528.037120117189 AND WFLWNS = 534.3208618164061 AND UALLEG3 = 326.5879156615625 AND ZNST81 = 98.369364874444 AND ZNST67 = 1.006589651107879 AND UALLEG = 152.5839221411094 AND ZNST72 = 88.75740202641016 AND ZNST78 = 0.01816449534961 AND ZNST89 = 1.054951740699205 AND UPRZ = 147.944380078125 AND QOVR = 1.08779780172795 AND UPRZ = 326.5836474669975 AND ZNST101 = 99.62434466936 AND KBCD015 = 15.0 AND ZNST79 = 100.0849699696795 AND KBCD07 = 20.0 AND BFPV488 = 0.924419710995938 AND WFLWNS = 534.562072739061 AND UALLEG = 100.5338652359061 AND BFPV479 = 0.01830245584204 AND QPRZH = 40.00765991109575 AND ZNST100 = 236.32865907617 AND ZNST74 = 65.308997744106 AND ZNST79 = 64.71767510986328 AND BFPV478 = 0.8725782036781311 AND ZNST1 = 3.33230302647981 AND QPRZLD = 0.9297760263212 AND UALLEG3 = 326.596822753906 AND ZVCT = 53.15314869869067 AND WCHGNO = 4.39840098040576 AND ZNST121 = 51.4104829259326 AND BFPV10 = 0.0 AND UALLEG1 = 390.4700927734575 AND ZNST15 = 0.230257815625 AND WDEWT = 8.91979694966455 AND UALLEG3 = 309.356567382125 AND BFPV499 = 0.87772276615428 AND UCOVD = 51.428163793946 AND ZNST79 = 66.9437879023363 AND UALLEG3 = 309.356567382125 AND BFPV489 = 0.43492591810173 AND ZNST77 = 50.03113174438477 AND ZNST75 = 65.308465860316 AND ZNST76 = 292.12677001953125 AND ZAPWTK = 9.71857929297363 AND UALLEG3 = 309.08474853156 AND BTV1418 = 0.0 AND ZNST48 = 0.294590471504883 AND FROEN = 60.0 AND PVAC = 721.2213134765625 AND KBCD022 = 99.0 AND BHTV2 = 1.0 AND PVCT = 1.733920469789425 AND KBCD011 = 160.0 AND ZNST36 = 29.253202078125 AND UALLEG3 = 309.3218078613281 AND BHTV6 = 1.0 AND ZNST3 = 0.0 AND KBCD019 = 1800.0 AND KBCD06 = 228.0 AND ZNST6 = 0.231696079964061 AND KBCD06 = 402.0 AND BHTV1 = 0.0 AND FSRMDPM = 1.68153817189780544 AND HCOIC = 0.0 AND BHTV22 = 0.0 AND BHTV30 = 0.0 AND BHTV13 = 0.0 AND KBCD08 = 228.0 AND ZNST22 = 2.0791811948342 AND BHTV = 1.0 AND KLAMPO15 = 1.0 AND BLV48 = 0.0 AND KBCD05 = 228.0 AND KBCD020 = 998.0 AND KBCD010 = 228.0 AND BTV143 = 1.0 AND KLAMPO48 = 0.0 AND KLAMPO9 = 0.0 AND KBCD04 = 228.0 AND KLAMPO21 = 1.0 AND KLAMPO24 = 1.0 AND BHTV1 = 0.0 AND KLAMPO19 = 0.0 AND KLAMPO18 = 0.0 AND KLAMPO17 = 0.0 AND KLAMPO28 = 1.0 AND KLAMPO29 = 0.0 AND KLAMPO = 0.0 AND KLAMPO31 = 1.0 AND KLAMPO16 = 1.0 AND KLAMPO15 = 0.0 AND KPV610 = 0.0
----- Subgroup 2 -----
Rule N° 1: IF BFPV145 = 0.652473628529656 AND ZCNDTR = 0.45296859741211 AND CIOADMPC = 99.29264831542969 AND UCHGUT = 211.9301453697188 AND BLV614 = 0.0 AND QPRZH = 140.129846434818645 AND ZREAC = 6.4624671936035165 AND ZNST85 = 534.489844373 AND ZNST79 = 66.7823394217811 AND UPRZ = 199.320313878956 AND ZVCT = 53.15314869869067 AND BHTV = 0.974524114942912 AND CAXOFF = 0.523381178616873 AND UPRZ = 340.19226591798875 AND CRETV = 1.0338860486993496407 AND ZNST85 = 534.287436136663 AND ZNST74 = 65.308997744106 AND ZNST87 = 528.1591798875 AND WSRFRAY = 0.0 AND BHTV = 1.0 AND BHRVCV = 0.99479681284324 AND ZNST66 = 0.43300340039063 AND ZNST103 = 87.0527252197266 AND ZNST86 = 534.317077657189 AND CIOADMPC = 2663.8357421875 AND UREDUT = 114.75112258594 AND ZCND = 0.781590044984846 AND ZNST108 = 39.89359614640234 AND ZNST68 = 1.006427407247097 AND ZNST1 = 99.2513107299048 AND BPRCV = 0.0 AND BHTV12 = 0.949208199977875
```

Fig. 4. A portion of the rule extraction results from the CNS data-based anomaly detection.

5. Conclusion

AI-based decision support systems are actively being researched to reduce human error in the operation of NPPs. However, the black-box characteristics of AI can be an obstacle to the practical application of such research. Recently, XAI methods have been applied to solve this problem. However, most XAI applications are limited to supervised learning-based algorithms. However, many decision researchers are also developing anomaly detection models using unsupervised learning algorithms. Therefore, applying XAI methods to unsupervised learning algorithms is necessary.

In this study, we consider the applicability of the OCSVM with rule extraction method, which combines the OCSVM method, an unsupervised learning algorithm, with the rule extraction method, an XAI algorithm. The case studies were performed to assess

the applicability of: (1) anomaly detection using accelerated aging data of IGBT and (2) anomaly detection using NPP simulation data. The extracted rules reveal the logical structure of the OCSVM method-based anomaly detection models. These results are expected to enhance the interpretability of unsupervised learning-based anomaly detection models. However, a limitation exists: each input variable contributes to a rule, potentially complicating interpretation as the number of variables increases. This limitation is difficult to address because we ultimately extract rules based on the characteristics of the OCSVM model. In future work, we plan to explore other AI models or supplement the rule reduction technique to better reflect the characteristics of the data and derive interpretable rules.

Acknowledgment

This work was supported by the Korea Institute of Energy Technology Evaluation and Planning (KETEP) grant funded by the Korea government (MOTIE) (20224B10100120, Development of commercialization technology for failure diagnosis of reactor control and digital I&C systems) and (20224B10100130, Development of operational state simulator for operating nuclear power plant and commercialization technology for artificial intelligence decision-making support system to prevent human error in accident operation).

REFERENCES

- [1] M. M. Moya, and D. R. Hush, Network Constraints and Multi-objective Optimization for One-class Classification, Neural Networks, Vol.9, pp.463-474, 1996.
- [2] A. Barbado, O. Corcho, and R. Benjamins, Rule Extraction in Unsupervised Anomaly Detection for Model Explainability: Application to OneClass SVM, Expert Systems with Applications, Vol.189, p. 116100, 2022.
- [3] J. Celaya, P. Wycocoki, and K. Goebel, IGBT Accelerated Aging Data Set, NASA Prognostics Data Repository, NASA Ames Research Center, Moffett Field, CA, 2009.