

Anomaly data synthesis and detection via domain randomization

Joonha Jun¹, and Jongsoo Lee¹

¹*School of Mechanical Engineering, Yonsei University, Seodaemun-gu, Seoul, 03722, Republic of Korea*

joonhajun@yonsei.ac.kr

jleej@yonsei.ac.kr

ABSTRACT

With the great development of artificial intelligence, the demand for a large amount of data necessary for learning is increasing. The synthesis of engineering data is challenging in that it is not only to combine data, but also to proceed with data synthesis while keeping the engineering characteristics intact. To address this problem, this paper proposes a synthesis and detection model of anomalous data utilizing domain randomization. This model learns data from existing systems to identify patterns via manipulated vectors and synthesizes new data by itself with domain randomization. With little amount of data, the trained model can accurately detect anomaly data in the system to apply prognostics and health management (PHM) in various environmental conditions for system's health analysis.

1. INTRODUCTION

Nowadays, using computer-aided design (CAD) tools and utilizing its simulation outcomes became important in engineering from product design to maintenance in the sense of economic efficiency. The need to reduce time and resources put in testing prototypes accelerated CAD development which led to early access to system's response data without demonstrating its actual physical form. In addition to development of CAD software, inconsistency between simulation and actual data has been gradually decreasing with artificial intelligence (AI) applications such as data augmentation or domain adaptation. However, to obtain high quality data output from AI models, sufficient amount of reliable actual data is needed to train them, which brings back the inefficient process of designing physical prototypes again. Therefore, algorithm that synthesizes data from small actual data and trains large domain of data itself is required.

In application of AI in the field of PHM, similar problems are frequently encountered since there lack sufficient anomaly data to train PHM algorithms. In this study, data

synthesis of time series data via domain randomization is developed to directly tackle the problem of data insufficiency. The algorithm is performed in the process of three steps. First, actual serial data is preprocessed with feature extraction and principal component analysis (PCA) to obtain PCA eigenvectors. Convolutional neural network (CNN) model is trained with eigenvectors to apprehend patterns of each data domain. Second, PCA eigenvectors are augmented with conditional generative adversarial network (CGAN) to fulfill the purpose of domain randomization. Lastly, generated eigenvectors are put in the PCA reconstruction converter to gain time series data type.

The rest of the sections are organized as follows. Section 2 describes the dataset of time series data and preprocessing step. Section 3, 4 shows PCA reconstruction and domain randomization used in the research respectively. Section 5 discusses how the application can be conducted and Section 6 concludes this research.

2. DATA PREPROCESSING

2.1. Dataset

The actual dataset is generated from testbed model GDMCAS-2000 made by DyLab corporation. The vibration signal is picked up from the bearings for about 2 minutes of 100Hz resolution. 9 types of data with difference in motor rpm with range of 300 to 1500. Time series data is generated and for preprocessing, length is manipulated to 2 seconds each. Figure 1 and Figure 2 shows the testbed used and an example of generated time series datum respectively.

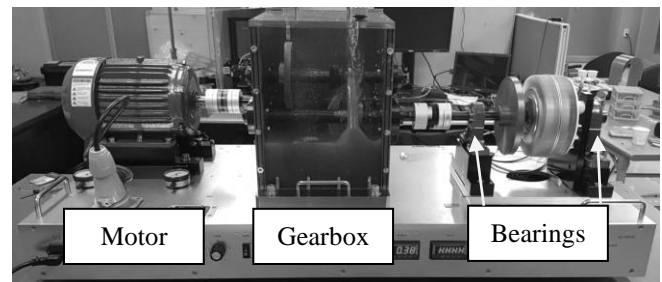


Figure 1. GDMCAS-2000 testbed

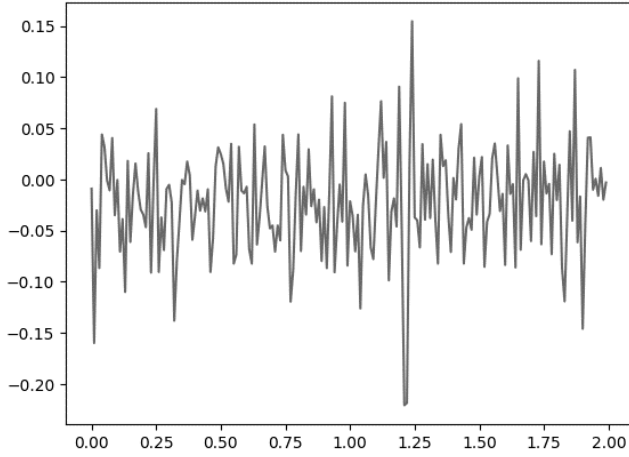


Figure 2. Time series datum example

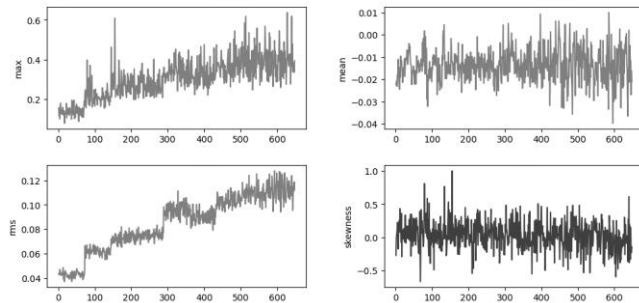
2.2. Feature Extraction

Feature extraction is mostly used well in machine learning where pattern recognition and image processing are done. In the case of where initial data is too large to be processed in the training, feature extraction is conducted to cut redundancy in the data and obtain reduced representation instead of the complete initial data.

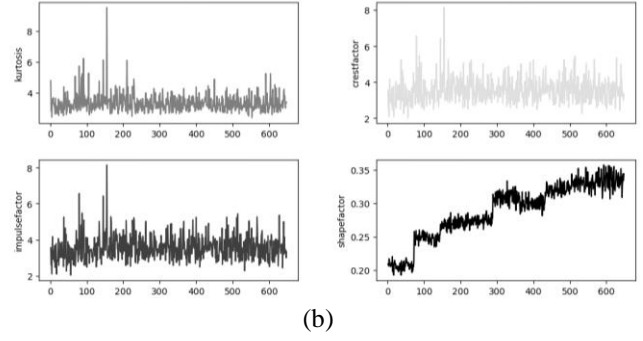
For this paper study, 144 seconds long 300, 500, 700, 800, 900, 1000, 1100, 1300, 1500 RPM data are used which are then divided into 2 seconds length each to create separate domains. Total of 648 data is utilized with such steps taken and each datum has length of 200 since data are acquired with 100Hz resolution.

To reduce dimension and manipulate data with PCA eigenvector, feature extraction must be preceded. 8 features of time series data are selected which consists of max, mean, root mean square, skewness, kurtosis, crest factor, impulse factor, and shape factor. Figure 3 shows feature extraction result of all 648 data in consideration.

From some features, data separation due to difference in RPM is evident but since not every external noise could not be controlled, it can be assumed from this stage that big noise has been picked up in the early stages of experiment.



(a)



(b)

Figure 3. Feature extraction of 648 time series data example (a) Max, mean, rms, and skewness features, (b) Kurtosis, crest factor, impulse factor, and shape factor features

2.3. Principal Component Analysis

Principal component analysis (PCA) combines features extracted to make new features, called principal components, that are orthogonal to each other. The first principal component contains the most features of the data, then the second, third, etc. Each datum will have an eigenvector after PCA process.

From time series data, feature extraction and PCA is conducted to make domain randomization possible with manipulating each PCA eigenvector and make large data domain with small amount of data. Figure 4 shows the correlation between 1st, 2nd, 3rd principal components respectively.

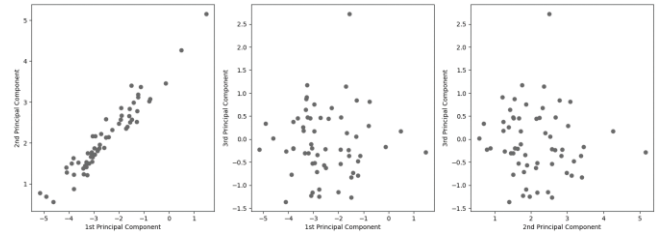


Figure 4. 1st, 2nd, 3rd principal component correlation results of 8 features of 648 time series data

3. PCA RECONSTRUCTION

Eigenvectors and covariance matrix results from PCA process contains most of the information of data. From this information, reconstruction is done by inverse operation of PCA formula. At this point, eigenvector is manipulated to make different output which will be used to train algorithms with anomaly label. Equation 1 refers to the reconstruction operation where R is reconstruction result, s is PC score, \mathbf{v}^T is transposed eigenvector, μ is mean.

$$R = s \cdot \mathbf{v}^T + \mu \quad (1)$$

For example, when using one RPM category of the data which numbers add up to 72 for each RPM, maximum 72 principal components can be extracted. With the use of 54

principal components, which has a coverage of 95% ($2\text{-}\sigma$ range), reconstruction is done as Figure 5.

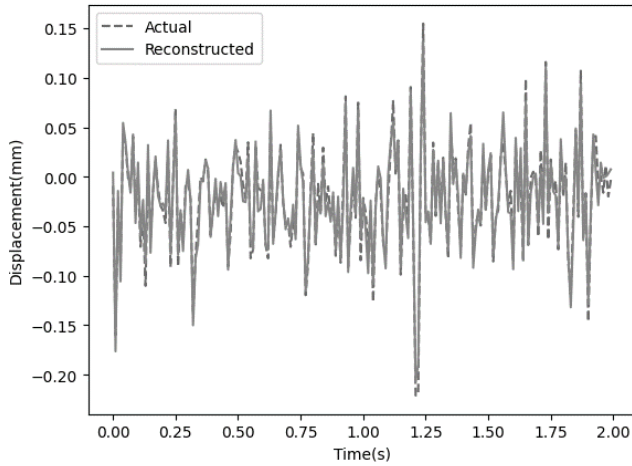


Figure 5. PCA reconstruction with 54 principal components

4. DOMAIN RANDOMIZATION

Domain randomization is a concept of expanding training data domain to improve accuracy of AI model which originates from paper by Josh Tobin et al. [1] In the paper, domain randomization took the part of generating many random situations under certain conditions so that the algorithm can train large domain of data without extra actual data.

In the same sense, domain randomization in this work operates in the purpose of data synthesis of similar patterns compared to actual data. With application of conditional generative adversarial network (CGAN) in augmenting and manipulating eigenvectors, domain randomization can be done with time serial data to robustly train anomaly detection classification models. Figure 6 shows the flow of domain randomization this research conducted.

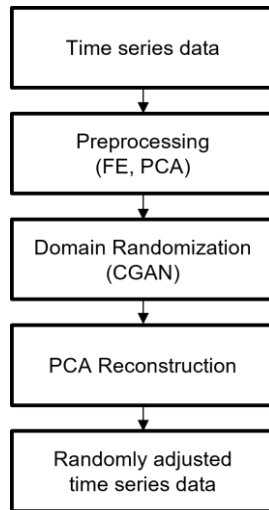


Figure 6. Flowchart of algorithm application

5. APPLICATION

With time series dataset and selected proper features, one can establish large and robust domain with small actual dataset via algorithm studied on this paper. As shown in Figure 6., pattern recognized time series data can be randomly augmented with certain guidance to use it alongside with actual data. Then, subtle differences made in actual data due to noise will affect less afterwards since those data will already be trained, leading to more accurate anomaly detection.

6. CONCLUSION

This paper utilized domain randomization via PCA and CGAN to efficiently apply it to time series data. The algorithm not only generates data but make large domain for robust training anomaly classification models. Method used for PCA eigenvector manipulation, in this study CGAN, can be varied to augment specific part of the data domain. From this study, more research on domain randomization on time series data is recommended to gain more insight on anomaly detection.

ACKNOWLEDGEMENT

This study was supported by the National Research Foundation of Korea (Grant No. 2022R1A2C2011034). This work was supported by Korea Evaluation Institute of Industrial Technology (Grant No. 20018208).

REFERENCES

- J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, P. Abbeel. (2017). Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World. *IEEE/RSJ International Conference on Intelligent Robots and Systems*.
- A. Pinceti, L. Sankar, O. Kosut. (2021), Synthetic Time-series Load Data via Conditional Generative Adversarial Networks. *IEEE Power & Energy Society General Meeting (PESGM)*.
- M. Ehrhart, B. Resch, C. Havas, D. Niederseer. (2022). A Conditional GAN for Generating Time Series Data for Stress Detection in Wearable Physiological Sensor Data. *Sensors*.
- Y. Chen, D. J. Kempton, A. Ahmadzadeh, R. A. Angryk. (2021). Towards Synthetic Multivariate Time Series Generation for Flare Forecasting. *ICAISC 2021, Artificial Intelligence and Soft Computing*.

BIOGRAPHIES

Joonha Jun in an Integrated M.S. and Ph.D. student in Mechanical Engineering at Yonsei University. His research interests are on the field of prognostics and health management (PHM), deep learning based data augmentation and synthesis via domain randomization.

Jongsoo Lee received his B.S. degrees in Mechanical Engineering from Yonsei University, Seoul, Korea, in 1988 and M.S. degree in Aerospace Engineering from the University of Minnesota, Minneapolis, MN, in 1992. He received Ph.D. in Mechanical Engineering from Rensselaer Polytechnic Institute, Troy, NY, in 1996. Dr. Lee is a

professor at School of Mechanical Engineering, Yonsei University, Seoul, Korea since Fall of 1997. His research areas include multi-physics design optimization (MDO) & prognostics and health management (PHM) with industrial artificial intelligence of data augmentation and knowledge transfer learning.