



**HAL**  
open science

# A Machine Learning Based Health Indicator Construction in Implementing Predictive Maintenance: A Real World Industrial Application from Manufacturing

Harshad Kurrewar, Ebru Turanoglu Bekar, Anders Skoogh, Per Nyqvist

► **To cite this version:**

Harshad Kurrewar, Ebru Turanoglu Bekar, Anders Skoogh, Per Nyqvist. A Machine Learning Based Health Indicator Construction in Implementing Predictive Maintenance: A Real World Industrial Application from Manufacturing. IFIP International Conference on Advances in Production Management Systems (APMS), Sep 2021, Nantes, France. pp.599-608, 10.1007/978-3-030-85906-0\_65 . hal-04022129

**HAL Id: hal-04022129**

<https://inria.hal.science/hal-04022129v1>

Submitted on 9 Mar 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



This document is the original author manuscript of a paper submitted to an IFIP conference proceedings or other IFIP publication by Springer Nature. As such, there may be some differences in the official published version of the paper. Such differences, if any, are usually due to reformatting during preparation for publication or minor corrections made by the author(s) during final proofreading of the publication manuscript.

# A Machine Learning Based Health Indicator Construction in Implementing Predictive Maintenance: A Real World Industrial Application from Manufacturing

Harshad Kurrewar, Ebru Turanoglu Bekar<sup>(✉)</sup>[0000-0003-4858-4386], Anders Skoogh<sup>[0000-0001-8519-0736]</sup>, and Per Nyqvist

Chalmers University of Technology, Hörsalsvägen 7A, 412 96 Gothenburg, Sweden  
ebrut@chalmers.se

**Abstract.** Predictive maintenance (PdM) using Machine learning (ML) is a top-rated business case with respect to the availability of data and potential business value for future sustainability and competitiveness in the manufacturing industry. However, applying ML within actual industrial practice of PdM is a complex and challenging task due to high dimensionality and lack of labeled data. To cope with this challenge, this paper presents a systematic framework based on an unsupervised ML approach by aiming to construct health indicators, which has a crucial impact on making the data meaningful and usable for monitoring machine performance (health) in PdM applications. The results are presented by using real-world industrial data coming from a manufacturing company. In conclusion, the designed health indicators can be used to monitor machine performance over time and further be used in a supervised setting for the purpose of prognostic like remaining useful life estimation in implementing PdM in the industry.

**Keywords:** Smart maintenance · Predictive maintenance · Health assessment · Machine learning · Feature selection and fusion · Real world industrial data

## 1 Introduction

Industrial digitalization is a key enabler of future competitiveness and sustainability in manufacturing companies. This allows for increased automation and data exchange through increased deployment of information and communication technologies (i.e., smart sensors) [1]. Recent studies have also shown that digital solutions in the maintenance field, specifically in Predictive Maintenance (PdM), have the highest potential to contribute towards industrial competitiveness with the vision of failure-free production [2]. Therefore, PdM research has a lot of attention in industry and academy due to its potential benefits in terms of reliability, safety and maintenance costs, among many other benefits and has been started to be adopted by many manufacturing companies [3]. Machine Learning

(ML), within Artificial Intelligence, has emerged as a powerful tool for developing intelligent predictive algorithms in many applications of manufacturing [4]. Therefore, ML provides powerful predictive approaches for implementing different PdM applications such as health indicators construction, (early) anomaly detection, and remaining useful life (RUL) estimation [5–7]. Although they have been successfully implemented to improve PdM capabilities, many of them follow a supervised learning approach. This perfectly works fine with experimental data sets, which means in the presence of a high amount of labelled data (i.e., the data is annotated with an actual machine health condition or the data contains examples of every possible fault type) [8]. However, the high amount of required labelled data for supervised prediction models in implementing PdM might not be available (partially/completely missing) in real-world industrial environments [9]. This shows that there still exists a gap, which most ML-based models for PdM are not designed with regard to actual industrial practice and are not validated with industrial data [9]. To overcome this gap, we propose a systematic framework to construct suitable health indicators that can be monitored and interpreted easily by practitioners in implementing PdM in real-world industrial environments. The main contribution of the proposed health indicator construction framework is that it gives significant insights into the implementation of unsupervised ML in PdM by identifying, cleaning, extracting, and selecting relevant features and fusing them as health indicators. Furthermore, it can also be effectively used with different condition monitoring data types with the extension of the framework in machine health assessment for PdM applications.

The remainder of the paper is structured according to the presented approach. In section 2, related literature focuses on health indicator construction in PdM. Section 3 presents the framework by describing the steps of the designed approach. Section 4 presents the results from real-world industrial application. Finally, section 5 concludes the paper with a summary and further research directions.

## 2 Related Literature

According to a systematic and comprehensive review done by [10], it was noted that ML algorithms have increasingly being applied for designing PdM applications especially handling the health status of industrial equipment and estimation of its RUL. Furthermore, a framework proposed by [11] for a data-driven based PdM generally covers three steps such as data acquisition and preprocessing, the construction of health indicator, and prognostics (i.e., prediction of a failure time). Among these steps, the construction of health indicators has a crucial impact on monitoring a system’s performance evolution over time. Moreover, a health indicator can be even used to not just detect the deviation from the healthy system conditions but also by trending it to predict the RUL of the system [12]. Moreover, there are two significant steps for a health indicator construction in literature: feature extraction and the selection and feature fusion, where multidimensional features are transformed into a one-dimensional health

indicator using ML algorithms [13]. For instance, [11] proposed a recurrent neural network (RNN) based health indicator to predict the RUL of bearings. The most sensitive features were selected from an original feature set based on monotonicity and correlation metrics and then fed into an RNN network to construct a health indicator, which was used RUL estimation of generator bearings from wind turbines. [14] also implemented RNN to fuse the selected sensitive features to construct a health indicator that incorporates mutual information of multiple features correlated with the damage and degradation of bearing. As one of the most popular unsupervised ML algorithm for dimension reduction, Principle component analysis (PCA) is often applied to fuse multiple features and construct health indicators for different systems such as bearings and cutting tools [13]. [15] recently presented a new approach for machine tool component health identification based on unsupervised ML. They trained different clustering models using time series feature extraction and demonstrated promising potential of the unsupervised techniques for machine tool supervision. Briefly, given the literature mentioned above, it is important to highlight that unsupervised ML approaches are particularly promising and need to be further developed for PdM applications in industries.

### 3 Methodology

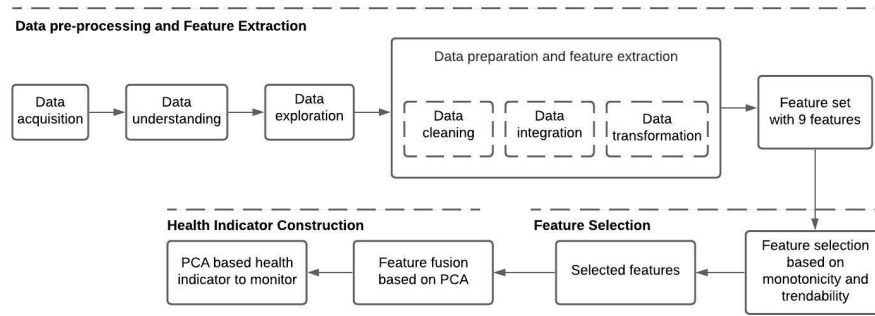
To establish a systematic framework for health indicator construction, the cross-industry standard process for data mining (CRISP-DM) is followed as a reference model in this paper. It provides a structured methodology for planning and managing the data-driven knowledge discovery process in data mining projects [16]. The CRISP-DM methodology consists of six iterative phases such as business understanding, data understanding, data preparation, modelling, evaluation and deployment. With these phases, the data is ensured regarding its relevance with identified business objectives in the business understanding phase of the process. However, it should be noted that the CRISP-DM methodology is taken a basis in this study. Hence, it is adjusted for practical development and implementation of the proposed framework according to real-world industrial application requirements. To answer this study's data-related questions, some certain phases of CRISP-DM, which are business understanding, data understanding, data preparation, and modelling, are focused and adapted to construct health indicators in the proposed framework. Therefore, the proposed framework can differ from the CRISP-DM with the main focus of the feature selection part. A flow chart of the proposed framework is illustrated in Figure 1.

Following understanding the industrial requirements (analysis goals) and exploring the current system and data, suitable data sources are chosen and pre-processed, including various data preprocessing tasks such as data cleaning and integration to get feature set. Afterwards, feature selection is performed to discard redundant features, which do not provide sufficient information related to the system performance (health) [11]. Therefore, the literature proposes some metrics such as monotonicity and trendability [13], which are used to select the

most important features from the feature set as shown in the proposed framework. The monotonicity metric evaluates an increasing or decreasing trend of the features, and it is calculated by the formula as follows [17].

$$Monotonicity = mean \left( \left| \frac{positive(diff(x_i)) - negative(diff(x_i))}{n - 1} \right| \right) \quad (1)$$

where n is measurement points with respect to time (cycles,  $i = 1, 2, \dots, n$ ) and diff is the difference of consecutive measurement points for each feature.



**Fig. 1.** A flow chart of the proposed framework for health indicator construction.

The trendability metric measures a linear correlation between features and operating time. It is calculated by the absolute minimum correlation in a population of the indicators using the formula as follows [17].

$$Trendability = \min(|\text{corrcoef}(x_i, x_j)|); i = 1, 2, \dots, m, j = 1, 2, \dots, m \quad (2)$$

where m is the number of features of the system.

These two feature metrics are restricted to the range  $[0, 1]$ , and they are positively associated with the feature performance accuracy, making them ideal as feature selection metrics. Therefore, the feature selection is performed based on a final feature importance score, which is computed as a linear combination of the above two metrics (considering each metric equally weighted) for each feature by using the formula as follows:

$$Final\_feature\_importance\_score = Monotonicity + Trendability \quad (3)$$

It should also be noted that the average value of these two metrics is computed for determining a threshold level to select the most important features.

In the final step of the proposed framework, the PCA is utilized to fuse the selected features for health indicator construction. PCA is a powerful unsupervised ML algorithm, which reduces the dimensionality of the data and retains

most of the variation (information) in the data set [18]. To take advantages of mutual information from the selected features, they are fed to the PCA algorithm to estimate the number of principal components that cover a minimum 85% cumulative percentage of explained variance and satisfy the "rule of thumb" inertial shift through scree plot visualization of eigenvalues. Based on this principle, the selected principal components are determined as health indicators to be monitored for machine performance and further be used in the estimation of RUL. As a common acknowledgement, the individual case study requires a different approach. This refers to that no single algorithm is superior to the other, as they each serve a purpose and are case-dependent. Despite this, the efforts in this study have been invested in developing a structured approach to create and streamline the model for health indicator construction. Furthermore, to approach the real-world case application effectively, and also with acknowledgment that the PCA is the most well know and applied technique in the literature [19], any comparisons between the PCA and the other types of unsupervised dimensionality reduction techniques have not been focused in this study.

## 4 A Real-World Industrial Application from Manufacturing

This section will give the results according to each step of the proposed framework for health indicator construction. Python programming language and its libraries such as Matplotlib and Seaborn for data visualization, Numpy and Pandas for data preprocessing and Sci-kit learn for machine learning are used to apply the proposed framework.

### 4.1 Data Preprocessing and Feature Extraction

The two bottleneck machines from an engine component line in one of the leading manufacturing industry in Sweden is used as a real-world industrial application in this paper. The industrial high challenging data (high frequency and lack of labels) coming from sensors and control system of the machines (machine programmable logic controllers (PLCs)) are collected by agents, converted into understandable structures and stored in a database system. This database system contains four structured query language (SQL) tables, including "Process Data" and "Vibration Data" from sensors, as well as machine-motor data from the control systems' including "Machine.1" and "Machine.2" tables. According to production and maintenance domain experts, vibration data should be used to generate health indicators because it is significant information when assessing the machine's condition practically. With this information given by the domain experts, the "Process Data" table is not considered for further analysis in this study.

We have vibration measurements such as Acceleration\_RMS, Velocity\_RMS, and Acceleration\_Peak, measured by different sensors on the spindle, spindle-motor, and gearbox of the machine, which is stored in the SQL table Vibration\_Data. These measurements have been collected at each 1-s interval for more

than one year. Machine-motor data is stored in the SQL databases of Machine.1 and Machine.2, including power consumption, torque drive, and motor temperature. These tables also contain some other values such as event\_time, vibration\_start\_stop, and block\_type as well. The scope of the data is given in Table 1. Here, Acceleration\_peak (\_A\_Peak) is a time-domain feature while Acceleration\_RMS (\_A\_RMS) and Velocity\_RMS (\_V\_RMS) are frequency domain features. The difference between the maximum and minimum values of the given time-domain based signals is defined as the Peak [20]. The root mean square (RMS) is used to calculate the signal’s overall energy and amplitude [20]. The data from PLCs are used for synchronization and preprocessing with the vibration data. Thereby, the vibration data is extracted at the beginning of each cycle, which means before the machining process starts by incorporating cycle start-stop information from PLCs. The domain experts suggest this period for analysis of the machine’s condition without any loading.

After relationships between the measurements given in Table 1 are explored by using visualization techniques and trend analysis based on the average smoothing method, it is observed that the spindle motor acceleration RMS value shows a slightly positive trend. In contrast, gearbox acceleration RMS shows a slightly negative trend in terms of time (cycles). This analysis motivated us to further evaluate each feature based on monotonicity and trendability metrics. It should be noted that the feature set is also normalized before comparing them with each other due to unit differences of the features by using a proper normalization method which is called z-score [21].

**Table 1.** Data Scope.

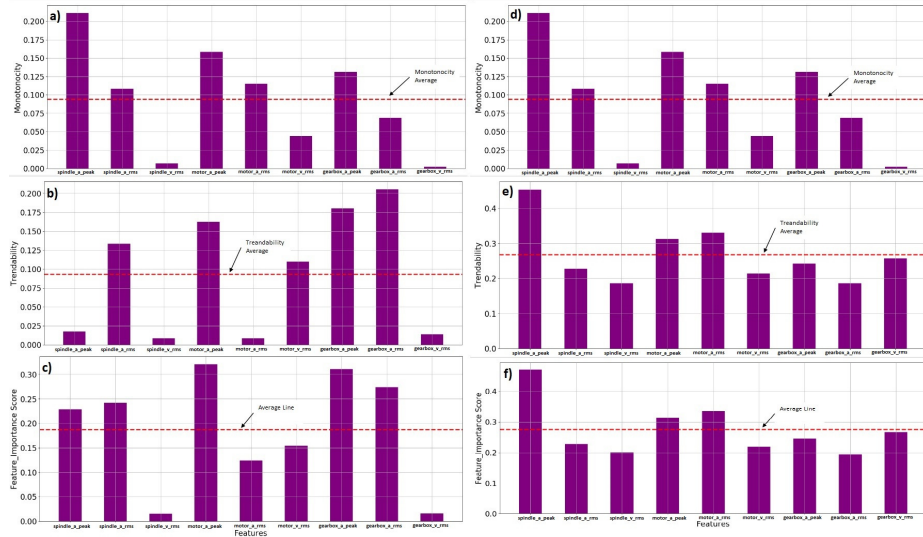
<b>Number of machines</b>	2
<b>Vibration measurements</b>	Spindle Spindlemotor Gearbox
<b>Computed features(time and frequency domain)</b>	Acceleration_RMS Velocity_RMS Acceleration_Peak
<b>Number of sensor measurements</b>	9
<b>Data resolution</b>	1-s
<b>Period of interest</b>	October 2019 - July 2020

## 4.2 Feature Selection

A good health indicator is characterized by monotonicity and trendability with respect to time (or cycles), as explained in the Methodology section of the paper. According to the final importance score computed from these two metrics, the features such as Spindle\_A\_Peak, Spindle\_A\_RMS, Spindlemotor\_A\_Peak, Spindlemotor\_A\_RMS and Gearbox\_A\_Peak are selected as the most important



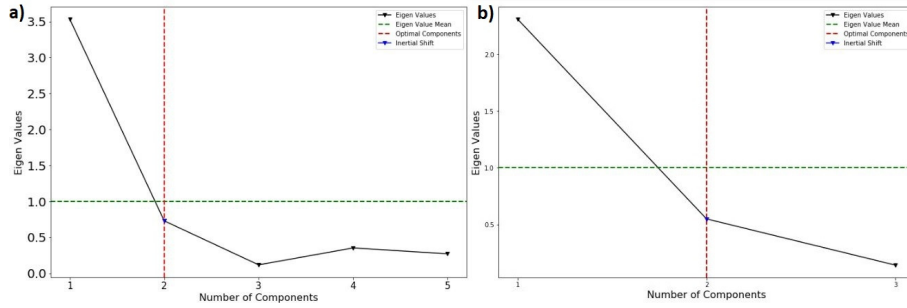
features for Machine\_1. For Machine\_2, Spindle\_A.Peak, Spindle\_A\_RMS and Spindlemotor\_A\_RMS are chosen as the most important features. The feature selection is graphically shown in Figure 2.



**Fig. 2.** (a),(d) Monotonocity (b),(e) Trendability and (c),(f) Final feature importance score of Machine\_1 and Machine\_2.

### 4.3 Health Indicator Construction

In this section, the PCA algorithm is used to fuse the selected features for both machines to construct health indicators. In this algorithm, the data is mapped linearly into a lower dimension space to maximize the variance of the data. The covariance matrix for the data is constructed, and the eigenvectors and eigenvalues of the covariance matrix are calculated. Then, based on the inertia shift of the principal component, feasible components are selected. The inertial shift of the principal components can be seen in the scree plot in Figure 3. The plot is calculated based on the eigenvalues on the ordinate and the number of components on the abscissa. Since there is an inertia shift after two principal components in the scree plot of Machine\_1 and Machine\_2 and the decay becomes more linear, two optimal components are selected, which explain more than 85% and 95% variation for Machine\_1 and Machine\_2, successively. Thus, the complex and high dimensional vibration data are simplified and scaled-down, two principal components to be used as health indicators for easily monitoring the performance of the machines over time (cycles). The robustness of the proposed framework is demonstrated in the results by identifying similar features in two



**Fig. 3.** (a) Scree plot of PCA for Machine\_1. (b) Scree plot of PCA for Machine\_2.

identical machines from the same production line. This reflects the degree of assurance that a case-specific health indicator construction may maintain its effectiveness and applicability in real-world industrial settings. In addition, the proposed framework advises and supports practitioners in the decision-making of machine health by visualizing the health indicators in an illustrated manner.

## 5 Conclusion

In this paper, a systematic framework is proposed to construct health indicator in implementing PdM. After data preprocessing and feature extraction based on industrial requirements, the most important (informative) features are selected based on their final feature importance score. These selected features are used as inputs into the PCA algorithm for designing suitable health indicators. We demonstrate our results from entirely real-world industrial data during the implementation of the proposed framework, which is a good contribution to fulfil the gap between research and actual industrial practice of PdM by dealing with the challenges such as high dimensionality and lack of labelled data. Therefore, this study would help practitioners in implementing PdM in real-world industrial environments. Further research is required to build predictive models for RUL estimation using the designed health indicators from this study.

**Acknowledgement.** The authors would like to thank the Production 2030 Strategic Innovation Program funded by VINNOVA for their funding of the research project SUMMIT - SUsustainability, sMArt Maintenance factory design Testbed (Grant No. 2017-04773), under which this research has been conducted. Thanks also to Anders Ramst orm and Robert Bergkvist, who supported real-time data from a real-world manufacturing system. This research has been conducted within the Sustainable Production Initiative and Production Area of Advance at the Chalmers University of Technology.

## References

1. Bokrantz, J., Skoogh, A., Berlin, C., Wuest, T., Stahre, J.: Smart maintenance: A research agenda for industrial maintenance management. *Int. J. Prod. Econ.*, **224**, 107547 (2020)
2. May, G., Kyriakoulis, N., Apostolou, K., Cho, S., Grevenitis, K., Kokkorikos, S., Milenkovic, J., Kiritsis, D.: Predictive Maintenance Platform Based on Integrated Strategies for Increased Operating Life of Factories. In: Moon, I., Lee, G.M., Park, J., Kiritsis, D., Cieminski, G.V. (eds.) *APMS 2018. IAICT*, vol. 536, pp. 279–287. Springer, Cham (2018). <https://doi.org/10.1007/978-3-319-99707-0>
3. Lee, J., Ni, J., Singh, J., Jiang, B., Azamfar, M., Feng, J.: Intelligent maintenance systems and predictive manufacturing. *ASME. J. Manuf. Sci. Eng.*, **142**(11), 1–23 (2020)
4. Wuest, T., Weimer, D., Irgens, C., Thoben, K. D.: Machine learning in manufacturing: advantages, challenges, and applications. *Prod. Manuf. Res.*, **4**(1), 23–45 (2016)
5. Carvalho, T. P., Soares, F. A., Vita, R., Francisco, R. D. P., Basto, J. P., Alcalá, S. G.: A systematic literature review of machine learning methods applied to predictive maintenance. *Comput. Ind. Eng.*, **137**, 106024 (2019)
6. Jimenez, J. J. M., Schwartz, S., Vingerhoeds, R., Grabot, B., Salaiin, M.: Towards multi-model approaches to predictive maintenance: A systematic literature survey on diagnostics and prognostics. *J. Manuf. Syst.*, **56**, 539–557 (2020)
7. Lughofer, E., Mouchaweh, S. M.: *Predictive maintenance in dynamic systems*. 1st edn. Springer, Switzerland (2019). <https://doi.org/10.1007/978-3-030-05645-2>
8. Olesen, J. F., Shaker, H. R.: Predictive maintenance for pump systems and thermal power plants: State-of-the-art review, trends and challenges. *Sensors*, **20**(8), 2425 (2020)
9. Zhai, S., Gehring, B., Reinhart, G.: Enabling predictive maintenance integrated production scheduling by operation-specific health prognostics with generative deep learning. *J. Manuf. Syst.*, (2021)
10. Çınar, Z. M., Nuhu, A. A., Zeeshan, Q., Korhan, O., Asmael, M., Safael, B.: Machine learning in predictive maintenance towards sustainable smart manufacturing in industry 4.0. *Sustainability*, **12**(19), 8211 (2020)
11. Guo, L., Li, N., Jia, F., Lei, Y., Lin, J.: A recurrent neural network based health indicator for remaining useful life prediction of bearings. *Neurocomputing*, **240**(31), 98–109 (2017)
12. Fink, O. : Data-driven intelligent predictive maintenance of industrial assets. In: Smith A. (eds) *Women in Industrial and Systems Engineering. Women in Engineering and Science*, pp. 589–605. Springer, Cham (2019). [https://doi.org/10.1007/978-3-030-11866-2\\_25](https://doi.org/10.1007/978-3-030-11866-2_25)
13. Lei, Y., Li, N., Guo, L., Li, N., Yan, T., Lin, J.: Machinery health prognostics: A systematic review from data acquisition to RUL prediction. *Mech. Syst. Sig. Process.*, **104**(1), 799–834 (2018)
14. Ning, Y., Wang, G., Yu, J., Jiang, H.: A feature selection algorithm based on variable correlation and time correlation for predicting remaining useful life of equipment using RNN. In: *Proceedings of the 2018 Condition Monitoring and Diagnosis (CMD)*, pp. 1–6. IEEE, Australia (2018)
15. Gittler, T., Scholze, S., Rupenyan, A., Wegener, K.: Machine tool component health identification with unsupervised learning. *J. Manuf. Mater. Process.*, **4**(3), 86 (2020)

16. Schröer, C., Kruse, F., Gómez, J. M.: A systematic literature review on applying CRISP-DM process model. *Procedia Comput. Sci.*, **181**, 526–534 (2021)
17. Saidi, L., Ali, J. B., Bechhoefer, E., Benbouzid, M.: Wind turbine high-speed shaft bearings health prognosis through a spectral Kurtosis-derived indices and SVR. *Appl. Acoust.*, **120**, 1–8 (2020)
18. Bekar, E. T., Nyqvist, P., Skoogh, A.: An intelligent approach for data pre-processing and analysis in predictive maintenance with an industrial case study. *Adv. Mech. Eng.*, **12**(5) (2020)
19. Tharwat, A.: Principal component analysis-a tutorial. *Int. J. Appl. Pattern Recognit.*, **3**(3), 197–240 (2016)
20. Atamuradov, V., Medjaher, K., Camci, F., Zerhouni, N., Dersin, P., Lamoureux, B.: Machine health indicator construction framework for failure diagnostics and prognostics. *J. Sign. Process. Syst.*, **92**, 591–609 (2020)
21. Jain, A., Nandakumar, K., Ross, A.: Score normalization in multimodal biometric systems. *Pattern Recognit.*, **38**(12), 2270–2285 (2005)