14th IMEKO TC10 Workshop Technical Diagnostics
New Perspectives in Measurements, Tools and Techniques
for system's reliability, maintainability and safety
Milan, Italy, June 27-28, 2016

# A Data-Driven Prognostic Approach Based on Sub-Fleet Knowledge Extraction

Giacomo Leone[1], Loredana Cristaldi[1], Simone Turrin[2]

[1] *Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB), Politecnico di Milano, Piazza Leonardo da Vinci, 32, 20133, Milano, Italy, {giacomo.leone, loredana.cristaldi}@polimi.it*
[2] *ABB AG, Corporate Research Center, Wallstadter Str. 59, 68526 Ladenburg, Germany, simone.turrin@de.abb.com*

*Abstract* – **In this paper, a data-driven prognostic algorithm for the estimation of the Remaining Useful Life (RUL) of a product is proposed. It is based on the acquisition and exploitation of run-to-failure data of homogeneous products, in the followings referred as fleet of products. The peculiar feature of the proposal is that the products composing the fleet are not strictly required to belong to the same system or plants, since the only constraint is that they are characterized by similar operating conditions (f.i. installed in the same region or operating in the same industrial application). The algorithm, indeed, is able to detect the set of products (sub-fleet of products) showing highest usage and degradation pattern similarity with the one under study and exploits the related monitoring data for a reliable prediction of the RUL, resulting in a potential tool for an effective Predictive Maintenance (PdM) strategy.**

## I. INTRODUCTION

The Remaining Useful Life (RUL) of a system is defined as the useful life left at a particular time instant that is the remaining time interval in which it will be able to meet its operating requirements. RUL estimation represents the core of the Prognostics and Health Management (PHM) programs which aim to a reduction of maintenance and life-cycle management costs, an increase of the systems availability and the adoption of Predictive Maintenance (PdM) strategies [1,2].

In the literature, the prognostic algorithms for the RUL estimation are usually classified in three different categories. The first class is related to the model-based approaches that refer to physical models describing the behavior of the systems under study. Such models can be very accurate but often require a strong and detailed knowledge of the inherent physics-of-failure. It follows that they are often very specific to the case study and their implementation is not always possible.

On the other hand, data-driven algorithms, the second main category found in the literature, are mainly based on the exploitation of the collected run-to-failure data and usually do not require particular knowledge about the inherent failure mechanisms. They provide a good trade-off between model complexity and results accuracy.

Finally, hybrid approaches attempt to leverage the advantages of combining the prognostics models in the aforementioned different categories for RUL prediction.

In the last years, data-driven approaches have experienced a wide diffusion. One reason is their suitability for applications related to complex engineered systems for which the definition of analytical models is a complex and resources demanding task. Another factor is the increasing availability of cheap monitoring systems that allow the collection of condition monitoring data in substantial quantities. A complete and detailed review about them is given in [3].

The same authors of this paper already presented two data-driven prognostic algorithms, one based on the statistical extraction and exploitation through Monte Carlo (MC) simulations of reliability and maintenance knowledge [4] and one based on a Machine Learning solution, in particular a Neural Network (NN) architecture [5]. A key concept and novel differentiator of the proposed approaches was the concept of fleet of products. A fleet of products is a set of homogenous products, with respect to the function for which they are intended, clustered together following different possible criteria such as belonging to the same customer, being installed in the same region or same industrial application and so on. The advantage of this practice is the possibility to extract fleet-specific usage and degradation profiles that can be exploited for the RUL prediction of a specific element (i.e. a specific product) of the selected fleet. In particular, the contribution of the acquisition of knowledge at fleet level on the improvement of the prognostic ability is quite relevant when the estimation of the RUL for a product at its early stage is of interest. Predicting the future behavior, in fact, is tied to the ability to learn from the past [1], and this is quite limited if the product is not in the mature stage of its life.

The application cases considered in [4-5], as well as in this paper, are Medium Voltage (MV) and High Voltage (HV) Circuit Breakers (CBs). MV and HV CBs are crucial elements in power transmission and distribution

systems. They are usually characterized by a long useful life and high reliability, but at the same time, even a single failure of them may cause severe damages, from technical, economical and safety point of view. It follows that reliable prognostic models are necessary for such kind of engineered systems. Defining physical models for such devices, however, is a time consuming and laborious task, since several different failure mechanisms depending on many parameters (number of completed opening operations, contacts degradation, short circuit current) may occur and also interfere with each other. All these aspects motivate the choice of applying data-driven algorithms for the estimation of their RUL.

As already said, the approaches proposed in [4-5] are based on the concept of fleet of products. The application case proposed, however, represents one of the many classes of products for which collecting a relevant number of run-to-failure curves is difficult. A striking example for this are the Vacuum Circuit Breakers (VCBs). They are based on a relatively recent technology and producers claim for them a Mean Time To Failure (MTTF) of over 30 years. It follows that acquisition of CM data spanning all the lifetime for many VCBs that can be considered to work in similar conditions and industrial applications is often not possible.

These considerations have driven the authors to propose in this paper a new approach that can be considered as an extension of the algorithm presented in [4]. The main contribution is the presentation of a method for the selection among a fleet of products of a subset of them, namely a sub-fleet of products, showing highest degradation pattern similarity with the product under study and for which the RUL estimation is required. Then, the related monitoring data are exploited for a reliable prediction of the RUL, offering a potential tool for an effective Predictive Maintenance strategy. The resulting advantage is that in this way the constraints for the definition of a fleet of products for prognostics aims become less strict so that the reference library can be constituted by run-to-failure data of products belonging to different customers, different geographical regions and industrial application. The proposed methodology, in fact, will automatically select the products with most affine degradation profiles, discarding the ones that would weight in negative way in the RUL estimation.

This article is structured in the following way: in Section II the proposed methodology for the selection of an appropriate sub-fleet is described and compared with some models found in the literature. In Section III the integration of the methodology in the existing framework of the approach presented in [4] is illustrated. The results obtained for the application case are reported and commented in Section IV, then the paper is ended with the conclusions.

## II.  SUB-FLEET SELECTION

The definition of a sub-fleet of products consists in the selection among a given set of products of a subset of them that show higher similarity, in terms of observed degradation in time, with respect to the item for which estimation of the RUL is required. In the literature, some methods for a sub-fleet definition already exist. In particular, they are based on a similarity-based approach that consists in the evaluation of the similarity between the test trajectory pattern (monitored degradation pattern for the item for which the RUL has to be predicted) and the reference trajectory patterns stored in the database and use the RULs of these latter to estimate the RUL of the former, accounting for how similar they are [6]. In [7], the authors propose the definition of a similarity coefficient based on the sum of the squared errors between the monitored test pattern and the reference trajectories. In particular, given a test product x and a reference item j, the similarity coefficient $s_{xj}$ is calculated as:

$$s_{xj} = \sum_{i=1}^{I}\left(h_{ji} - h_{xi}\right)^2 \qquad (1)$$

where $h_{xi}$ ($h_{ji}$) is the observed degradation for the test product $x$ (library specimen $j$) at cycle $i$ and $I$ is the number of observed cycles. In the estimation of the RUL, in order to give more weight to the library specimens with larger similarity coefficients, a weighting coefficient for a product $j$ is defined as:

$$d_{xj} = \exp\left(-\frac{s_{xj}}{\alpha}\right) \qquad (2)$$

where $\alpha$ is selected according to the desired selectivity ( i.e. if $\alpha$ is small, few specimens are influential). Finally, the *RUL* for the product $x$ is computed according to (3):

$$RUL_x = \frac{\sum_{j=1}^{J} d_{xj} RUL_j}{\sum_{j=1}^{J} d_{xj}} \qquad (3)$$

This approach has been used also in [8], whereas in [9] a slight modification is proposed. The modifications are made in the RUL estimation (i.e. Eq. (3)), in which the most similar K percent number of the library samples are utilized rather than using whole dataset.

A different approach is suggested in [10]. The authors, in fact, propose a definition of a deterministic model $M_i$ for each i-th training unit of the library, so that, for a given time $t$, an estimated value $y$ of the Health Indicator (HI) variable that describes the degradation pattern of the item is provided.

At this point, if a sequence $Y = y_1, y_2, \ldots, y_r$ of values of the HI for a test unit is available, a distance measure between the model $M_i$ and Y is defined as the sum of the squared errors between the monitored test pattern and

estimations provided by the model, divided by the prediction variance of the model itself. Then the RUL estimation for the test product is equal to a weighted sum of the RUL of the reference products. The weights can be assigned according to different principles. One of them is to apply the k-nearest neighbor method that is to select the products with the k smallest distance values and apply a weight 1/k to their RULs.

In our proposal, the selection of the sub-fleet is carried out according to the distance concept, analogously to the previous discussed papers. In particular, the distance is computed comparing the Health Condition (HC) profiles over time of the test product and the fleet products. Such profiles are obtained from the Condition Monitoring (CM) data that provide direct or indirect information on the HC of the related products. For the application cases considered in this paper (i.e. medium voltage and high voltage circuit breakers) examples of CM data are measurement of the contact ablation, $SF_6$ gas density for gas insulated circuit breakers, and temperature of the interrupting chamber. In the followings, this notation will be used: HC=100% means that the product is in perfect healthy state, whereas HC=0% means that it has reached its End of Life (EoL), referring to a scenario in which it is not anymore able to perform its intended function, so that a maintenance activity, refurbishment, replacement or its disposal is required.

Let us now suppose that for the test product for which the RUL estimation is wished a partial monitoring of its HC profile is available and it is constituted by $n$ consecutive observations. The knowledge about the degradation profile of such product can be described through a time series $K_x$ composed by $n$ couples of values as follows:

$$K_x = \left\{ (t_1, HC_1), (t_2, HC_2), ..., (t_n, HC_n) \right\} \qquad (4)$$

In particular, the generic value $t_i$ represent the life stage (time instant or number of cycles) of the product $x$ at which the i-th observation about its HC has been carried out, being $HC_i$ the related value. At this point, the HC profile of the generic j-th fleet product can be compared with the test profile, determining the related time stamps corresponding to the HC values available for the test profile, namely $HC_1$, $HC_2$,…, $HC_n$. In other words, for each j-th product is possible to determine a time series $K_j$ defined as:

$$K_j = \left\{ (t_{1,j}, HC_1), (t_{2,j}, HC_2), ..., (t_{n,j}, HC_n) \right\} \qquad (5)$$

It is easy to understand that the more the HC profiles of the reference product $j$ and the test product $x$ are similar, the smaller is the difference between corresponding time instants, such as $t_1$ and $t_{1,j}$, $t_2$ and $t_{2,j}$ and so on. Starting from this assumption a distance value $d_{xj}$ that correlates the two products can be computed as:

$$d_{xj} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} \left( t_i - t_{i,j} \right)^2} \qquad (6)$$

that is the Root Mean Square Error (RMSE) between the two HC profiles. A small $d_{xj}$ means that the two profiles are similar or, equivalently, similar degradation processes characterize the two products, whereas large distance values are related to products that are subject to different degradation mechanisms and that should be excluded in the estimation of the RUL for the target products since are representative of different working conditions. At this point if the reference fleet is composed by $N_f$ products, it is possible to define a variable $d_f$, namely the fleet distance, that describes the distribution of the distance values for the fleet. The experimental sample for such variable is provided by the $N_f$ different values of the distance computed according to (6):

$$d_f = \left\{ d_{x1}, d_{x2}, ..., d_{xN_f} \right\} \qquad (7)$$

In order to isolate the sub-fleet of products showing highest functioning similarity with the one under study and exploit the related monitoring data for a reliable prediction of the RUL, a threshold $d_\alpha$, corresponding to the $\alpha$-quantile for the variable $d_f$ can be defined. The parameter $\alpha$ is substantially the desired degree of selectivity (i.e. the smaller the parameter $\alpha$, the stricter the selection criterion). Finally, the products characterized by a distance value smaller than $d_\alpha$ will compose the desired sub-fleet.

A crucial difference between our proposal and the solutions presented in the literature is that after the selection of the most representative sub-fleet for the target product, its RUL estimation is not just defined as a weighted sum of the RUL of such subset of products, but, on the contrary, a prognostic model that exploits all the information enclosed in their HC profile is involved. Taking into account the entire HC profiles instead of considering exclusively the RUL values, which represent only the last points of such curves, enables the possibility to include in the prognostic model a more complete information, so that, besides the RUL, also other outcomes of interest such as the Probability of Failure (PoF) within a predetermined window of time can be provided. Another benefit deriving from the application of this approach is that the sub-fleet is not strictly required to be solely composed by products with a known RUL (i.e. already failed), but also products characterized by a partial HC profile knowledge can be included in the analysis, extending the potential set of products from which extracting information suitable for prognostics issues. This factor is as more determining as more difficult is the acquisition of run-to-failure data in substantial quantities (f.i. this is the case for VCBs). In the next Section III, the proposed prognostic algorithm based on MC simulations is presented. A more detailed

description, however, is given in [4].

## III. PROGNOSTIC ALGORITHM DESCRIPTION

The main assumption underlying the algorithm is that the future usage of the test product (for which a partial observation of its HC profile is available) might possibly be similar to the usage profile of the sub-fleet of homogeneous/similar products. It follows that the first step for the enhancing of the precision in forecasting the target RUL is to extract knowledge from the condition monitoring data of such products. It has been already demonstrated in [4] that this approach allows to overcome the risk of carry out a long-term prediction relying on the limited portion of CM data that would be available if only the past history of the test product was considered. Finally, this knowledge can be exploited in order to predict the future HC profile over time and extracting a confidence interval for the test product RUL.

### A. Knowledge extraction at sub-fleet level

Let us suppose that after the sub-fleet selection described in Section II the historical data analysis is limited to $N_{sf}$ products (the subscript "sf" stays for "sub-fleet"). One way extraction to extract past usage information about such products is to take into account for each of them the distribution of the sampling time and the distribution of the health variation in the related CM data. The sampling time corresponds to the time interval between two consecutive observations, whereas the health variation is the variation observed in the HC profile during the same time interval. Considering the generic j-th product of the sub-fleet, the related sampling time vector is expressed as:

$$\overline{\Delta t_j} = \left\{ t_{2,j} - t_{1,j}, ..., t_{n_j,j} - t_{n_j-1,j} \right\} \qquad (8)$$

where $n_j$ is the number of available observation for that product. Similarly, the vector of health condition variation is expressed as:

$$\overline{\Delta HC_j} = \left\{ HC_{2,j} - HC_{1,j}, ..., HC_{n_j,j} - HC_{n_j-1,j} \right\} \qquad (9)$$

Considering MV and HV CBs, the variation of the sampling time provides information on the usage profile of the breaker which is equivalent to the frequency of the opening operations. On the other hand, the variation of the health condition represents the degradation profile of the product.

The extraction of sampling time and health variation information at sub-fleet level is carried out in systematic way applying (8) and (9) to each product of the selected subset, so that the following vectors are obtained:

$$\overline{\Delta t_{sf}} = \left\{ \overline{\Delta t_1}, \overline{\Delta t_2}, ..., \overline{\Delta t_{N_{sf}}} \right\} \qquad (10)$$

$$\overline{\Delta HC_{sf}} = \left\{ \overline{\Delta HC_1}, \overline{\Delta HC_2}, ..., \overline{\Delta HC_{N_{sf}}} \right\} \qquad (11)$$

as concatenation of the vectors obtained at single product level. Starting from the above vectors, the Cumulative Distribution Functions (CDFs) of the variation of the sampling time and health condition at sub-fleet product level is obtained.

### B. Knowledge exploitation

In order to exploit the extracted knowledge an approach based on Monte Carlo simulations is discussed. The procedure is explained as follows:
1. Generation of two random numbers $r_1$ and $r_2$, drawn from a uniform distribution with values between 0 and 1.
2. The two extracted numbers are used to determine the subsequent point in the HC vs. time curve for the target product. In particular, the next sampling time $\Delta t^*$ and health condition variation $\Delta HC^*$ are determined exploiting the related CDFs, applying the Inverse Transform Method [11]. The generation of two uniformly distributed random numbers is made starting from the assumption that the sampling time and the variation of the health condition are not correlated. Step 1 and 2, however, can be easily modified to take into account the case in which these two variables are correlated.
3. Steps 1 and 2 are iteratively repeated until the estimation of the HC reaches the value 0%. The corresponding estimated time provides an estimation of the EoL for the target product. The distance between the end of the observation window and the estimated EoL corresponds to an estimation of the test product RUL.
4. Steps 1 to 3 are repeated for a statistically significant number of times $N_{MC}$.

At the end of step 4, $N_{MC}$ estimations of the product RUL are available, so that a confidence interval at a desired confidence level can be obtained.

## IV. ALGORITHM VALIDATION

The application case for the presented approach are MV and HV CBs. The data used and reported in the contribution are confidential information (ABB property), hence the exact numerical values are not reported.

A set of CM data related to a fleet of 90 products coming from different customers, operating regions and applications is considered. The algorithms have been tested according to the following procedure:
1. Set one CB as test product
2. Set as reference fleet a number $N_f$ of products chose randomly among the original fleet.
3. Delimiting the observation window for the test product at a given percentage of its actual lifetime.
4. Set a value $\alpha$ and extracting the most affine sub-

fleet products among the fleet defined at the step 2.
5. Run the prognostic algorithm described in Section III and obtain a confidence interval for the RUL.
6. Repeat steps 4-5 varying $\alpha$ from 0.05 to 0.7, with incremental steps of 0.05.
7. Repeat steps 3-6 varying the observation window from 10% to 90% of the actual lifetime, with incremental steps of 10%.
8. Repeat steps 2-7 10 times extracting randomly another reference fleet of $N_f$ products. This is done in order to simulate different possible scenarios.
9. Repeat steps 2-8 varying $N_f$ from 9 to 81, with incrementing steps of 12.

The procedure just proposed is repeated setting cyclically as test product a different CB. The objective of the study is to analyze the improvements obtained in the estimation of the RUL for a given target product through the definition of a sub-fleet of products among a wider reference set. In doing this, particular attention is given to the effect of the parameter $\alpha$ that defines the desired degree of selectivity and the number of products that constitute the starting fleet. In Fig.1, below, the results obtained defining reference fleets of 21 products are depicted. The focus is on the variation of the results as function of the observation window length (the percentage of collected data for the test product with respect to its actual lifetime) and the set value of $\alpha$. In particular, with algorithm performance for a given observation window length and value of $\alpha$ is intended the percentage of products for which a correct estimation of the RUL (i.e. the estimated confidence interval of the RUL encloses the actual RUL value) is provided. The results refer to estimated confidence interval for the RULs at 95% level.
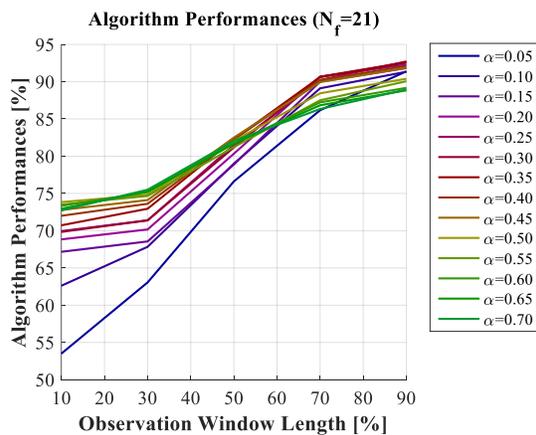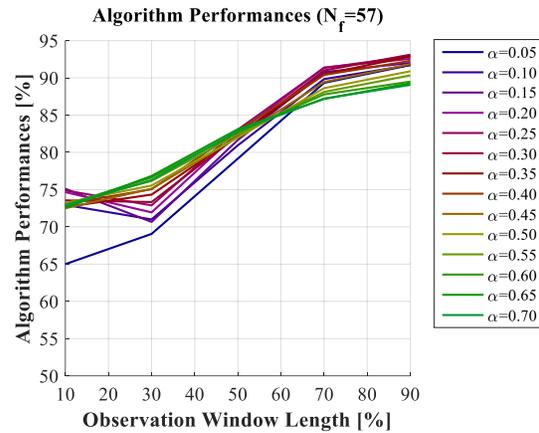


Fig. 2 Algorithm performances with $N_f$=57 and varying $\alpha$

test product) the best results are obtained with high values of $\alpha$, that means choosing a low selectivity criterion. On the contrary, as the observation window length increases, the best choice is to set lower values of $\alpha$. A motivation for this is the fact that as the amount of CM data increases the differences between the HC profiles of the test product and reference products are more emphasized, so that setting a low value $\alpha$ (high selectivity) allows to exclude effectively from the analysis the products showing a different degradation profile over time. This is not the case when few monitoring data are available. In such scenario a small portion of the HC-time space is explored so that often is not possible to get a solid evidence about different usage and degradation profiles for the products. It follows that being too selective would exclude from the analysis products that actually have HC profiles similar to that of the future test product, resulting underperforming.



Fig. 1 Algorithm performances with $N_f$=21 and varying $\alpha$

In an analogous way, Fig. 2 shows the results obtained with $N_f$=57. It is possible to observe that in both cases ($N_f$=21 and $N_f$=57), when a reduced observation window is considered (i.e. few available monitoring data for the
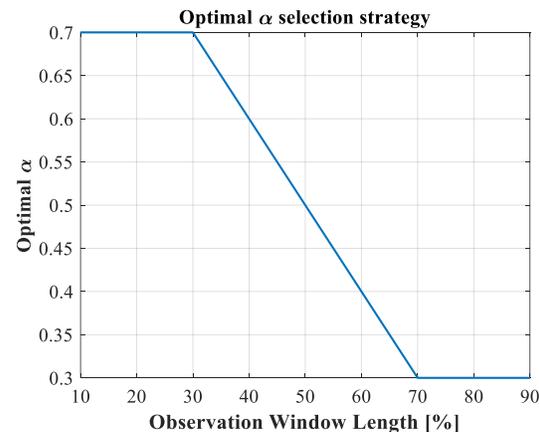


Fig. 3 Proposed rule of thumb for an optimal choice of $\alpha$

In Fig.3 it is proposed a rule of thumb for an optimal choice of the parameter $\alpha$ as function of only the

observation window length, no matter the number of products composing the reference fleet.

These results have been obtained from the analysis of the performances exhibited by the algorithm in the different configurations and offers a good trade-off between the simplicity to be applied and the resulting algorithm prognostic power.

Finally, Fig.4 compares the results about the algorithm performances obtained with $N_f$=81 in two opposing conditions. The black line depicts the results obtained considering a reference fleet of 81 products and applying the algorithm proposed in Section III neglecting the preliminary phase of sub-fleet discrimination. Conversely, the red line is related to the performances shown by the prognostic algorithm when it is run after the sub-fleet selection. In particular, such phase has been carried out following the strategy proposed in Fig.3. The results highlight how the proposed method to discriminate, among a reference fleet of products, a subset of them showing a closer degradation profile with respect to the product under test is actually very effective and contributes to the improvement of the prognostic power of the algorithm.

Another important result highlighted in the previous figures is that the sub-fleet selection methodology allows to get similar algorithm performances even starting from reference fleets composed by a quite different number of products.
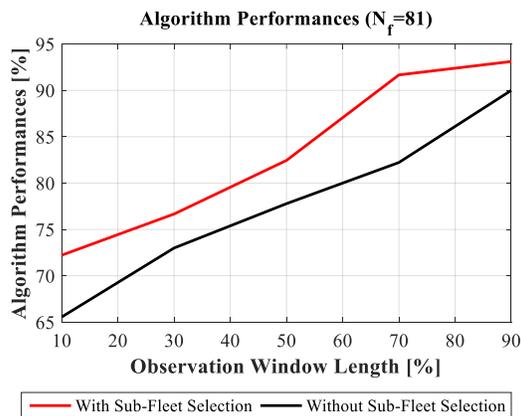


*Fig. 4 Improvement of algorithm prognostic power through a proper sub-fleet selection*

## V. CONCLUSIONS

In this paper, a data-driven prognostic algorithm for the estimation of the RUL of a product is proposed. It is based on the acquisition and exploitation of run-to-failure data of homogeneous products, referred as fleet of products. The core of the contribution is the proposal of a method for the identification of the products (sub-fleet of products) showing highest usage and degradation profile similarity over time with the one under study. The results

obtained for the application case of Medium Voltage and High Voltage Circuit Breakers have demonstrated that this approach contributes to the improvement of the RUL predictions and the increase of the prognostic algorithm performances. The methodology has been compared with those found in literature and the main differences have been highlighted. In particular, the advantages of our proposal is the possibility to exploit in the estimation of the future degradation of the test product all the information stored in the Health Condition profile of the sub-fleet products and not only the knowledge of their RUL, that corresponds to the last point of such curves. This makes the prognostic algorithm able to provide, besides the RUL, also other outcomes of interest such as the Probability of Failure (PoF) within a predetermined window of time. Another arising advantage is the possibility to involve in the analysis also reference products that still have not failed, making the approach very interesting for those classes of systems, such as Vacuum Circuit Breakers, for which the acquisition of run-to-failure data in substantial quantities is a hard task.

## REFERENCES

[1]    M. G. Pecht, Prognostics and Health Management of Electronics: *John Wiley & Sons*, 2008.

[2]    D. C. Swanson, "A general prognostic tracking algorithm for predictive maintenance," in *2001 IEEE Aerospace Conference Proceedings*, pp. 2971–2977.

[3]    M. Schwabacher, "A survey of data-driven prognostics," *Proc. AIAA Infotech@Aerosp. Conf*, Reston, VA, 2005.

[4]    S. Turrin, S. Subbiah, G. Leone, and L. Cristaldi, "An algorithm for data-driven prognostics based on statistical analysis of condition monitoring data on a fleet level," in *2015 IEEE International Instrumentation and Measurement Technology Conference (I2MTC)*, pp. 629–634.

[5]    L. Cristaldi, G. Leone, R. Ottoboni, S. Subbiah, S. Turrin, "A comparative Study on Data-Driven Prognostic Approaches Using Fleet Knowledge," in *2016 IEEE International Instrumentation and Measurement Technology Conference (I2MTC)*, pp. 263–268.

[6]    L. Angstenberger, Dynamic fuzzy pattern recognition with applications to finance and engineering. *Boston: Kluwer Academic, 2001*.

[7]    E. Zio, F. Di Maio, "A Data-Driven Fuzzy Approach for Predicting the Remaining Useful Life in Dynamic Failure Scenarios of a Nuclear Power Plant", *Reliability Engineering and System Safety, RESS*, 10.1016/j.ress.2009.08.001, 2009.

[8]    B. K. Guépié and S. Lecoeuche, "Similarity-based residual useful life prediction for partially unknown cycle varying degradation," in *2015 IEEE Conference on Prognostics and Health Management (PHM)*, pp. 1–7.

[9]    Eker, O.F., Camci, F. and Jennions, I.K., 'A Similarity-based Prognostics Approach for Remaining Useful Life Prediction', *The 2nd European Conference of the Prognostics and Health Management (PHM) Society*, Nantes, France, 8-10 July 2014, vol. 5, no. 11, 2014.

[10]    T. Wang, Jianbo Yu, D. Siegel, and J. Lee, "A similarity-based prognostics approach for Remaining Useful Life estimation of engineered systems," in *2008 International Conference on Prognostics and Health Management (PHM)*, pp. 1–6.

[11]    L. A. Grzelak, J. Witteveen, M. Suarez-Taboada, and C. W. Oosterlee, "The Stochastic Collocation Monte Carlo Sampler: Highly Efficient Sampling from 'Expensive' Distributions," *SSRN Journal*, 2014.