

Review and Analysis of Algorithmic Approaches Developed for Prognostics on CMAPSS Dataset

Emmanuel Ramasso¹ and Abhinav Saxena²

¹ *FEMTO-ST Institute, Dep. AS2M/DMA, UMR CNRS 6174 - UFC / ENSMM / UTBM, 25000 Besançon, France
emmanuel.ramasso@femto-st.fr*

² *SGT Inc., NASA Ames Research Center, Intelligent Systems Division, Moffett Field, CA, 94035-1000, USA
abhinav.saxena@nasa.gov*

ABSTRACT

Benchmarking of prognostic algorithms has been challenging due to limited availability of common datasets suitable for prognostics. In an attempt to alleviate this problem several benchmarking datasets have been collected by NASA's prognostic center of excellence and made available to the Prognostics and Health Management (PHM) community to allow evaluation and comparison of prognostics algorithms. Among those datasets are five C-MAPSS datasets that have been extremely popular due to their unique characteristics making them suitable for prognostics. The C-MAPSS datasets pose several challenges that have been tackled by different methods in the PHM literature. In particular, management of high variability due to sensor noise, effects of operating conditions, and presence of multiple simultaneous fault modes are some factors that have great impact on the generalization capabilities of prognostics algorithms. More than 70 publications have used the C-MAPSS datasets for developing data-driven prognostic algorithms. The C-MAPSS datasets are also shown to be well-suited for development of new machine learning and pattern recognition tools for several key preprocessing steps such as feature extraction and selection, failure mode assessment, operating conditions assessment, health status estimation, uncertainty management, and prognostics performance evaluation. This paper summarizes a comprehensive literature review of publications using C-MAPSS datasets and provides guidelines and references to further usage of these datasets in a manner that allows clear and consistent comparison between different approaches.

1. INTRODUCTION

In the past decade the science of prognostics has fairly matured and the general understanding of health prediction prob-

Emmanuel Ramasso et al. This is an open-access article distributed under the terms of the Creative Commons Attribution 3.0 United States License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

lem and its applications has greatly improved. Both data-driven and physics based methods have been shown to possess unique advantages that are specific to application contexts. However, until very recently, a common bottleneck in development of data-driven methods was the lack of availability of run-to-failure data sets. In most cases real-world data contain fault signatures for a growing fault at various severity levels but no or little data capture fault evolution all the way through failure. Procuring actual system fault progression data is typically time consuming and expensive. Fielded systems are, most of the time, not properly instrumented for collection of relevant data or are unable to distribute such data due to proprietary constraints. The lack of common data sets, which researchers can use to compare their approaches, has been an impediment to progress in the field of prognostics. To tackle this problem the Prognostics Center of Excellence (PCoE) at NASA's Ames Research Center established a prognostics data repository back in 2007 (Saxena & Goebel, 2008). Several datasets have been since published that have been used by researchers around the world. Among these datasets are five datasets from a turbofan engine simulation model - C-MAPSS (Commercial Modular Aero-Propulsion System Simulation) (Frederick, DeCastro, & Litt, 2007). By simulating a variety of operational conditions and injecting faults of varying degree degradation datasets were generated for prognostics development (Saxena, Goebel, Simon, & Eklund, 2008a). One of the first datasets was used for a prognostics data challenge at the PHM'08 conference. A subsequent set was then released later with varying degrees of complexity. These datasets have since been used very widely in publications for benchmarking prognostics algorithms.

The turbofan degradation datasets have received over seven thousand unique downloads in the last five years but algorithms developed using these have been published in only about seventy publications. Furthermore, in many publications it is not clear how authors are computing results and comparing with others. There has been a confusion and in-

consistency in how these datasets have been interpreted and used in many cases. Consequently, not all comparisons of performance can be considered valid. Therefore, this paper intends to analyze various approaches that researchers have taken to implement prognostics using these turbofan datasets. Some unique characteristics of these datasets are also identified that led to use of certain methods more often than others. Specifically, various differences among these datasets are pointed out. A commentary is provided on how these approaches fared compared to the winners of the data challenge. Furthermore, this paper also attempts to clear several issues so researchers in the future can take these factors into account in comparing their approaches with the benchmarks.

2. C-MAPSS DATASETS

C-MAPSS is a tool, coded in the MATLAB-Simulink[®] environment for simulating engine model of the 90,000 lb thrust class (Frederick et al., 2007). Using a number of editable input parameters it is possible to specify operational profile, closed-loop controllers, environmental conditions (various altitudes and temperatures), etc. Additionally, there are provisions to modify some efficiency parameters to simulate various degradations in different sections of the engine system.

2.1. Datasets characteristics

Using this simulation environment five datasets were generated. By creating a custom code wrapper, as described in (Saxena, Goebel, et al., 2008a), selected fault injection parameters were varied to simulate continuous degradation trends. Data from various parts of the system were collected to record effects of degradations on sensor measurements and provide time series exhibiting degradation behaviors in multiple units. These datasets possess unique characteristics that make them very useful and suitable for developing prognostic algorithms.

1. Data represent a multi-dimensional response from a complex non-linear system from a high fidelity simulation that very closely models a real system.
2. These simulations incorporated high levels of noise introduced at various stages to accommodate the nature of variability generally encountered.
3. The effects of faults are masked due to operational conditions, which is yet another common trait of most operational systems.
4. Data from plenty of units is provided to allow algorithms to extract trends and build associations for learning system behavior useful for predicting RULs.

Without a doubt these datasets were geared towards data-driven approaches where very little or no system information was made available to PHM developers.

As described in detail in Section 3, the analysis on the publications using these datasets shows that many researchers have tried to make comparisons between results obtained from these similar yet different datasets. This section briefly describes and distinguishes the five datasets and explains why it may or may not be appropriate to make such comparisons. Table 1 summarizes the five datasets. The fundamental difference between these datasets is attributed to the number of simultaneous fault modes and the operational conditions simulated in these experiments. Datasets #1 through #4 incorporate an increasing level of complexity and may be used to incrementally learn the effects of faults and operational conditions. Furthermore, what sets these four datasets apart from the challenge datasets is the availability of ground truth to measure performance. Datasets 1 – 4 consist of a *training set* that users can use to train their algorithms and a *test set* to test the algorithms. The ground truth RUL values for the test set are also given to assess prediction errors and compute any metrics for comparison purposes. Results between these datasets may not always be comparable as these data simulate different levels of complexity, unless a universal generalized model is available that regards datasets 1 – 3 as special cases of dataset #4.

The PHM challenge datasets are designed in a slightly different way and divided into three parts. Dataset #5T contains a *train set* and *test set* just like for datasets 1 – 4 except with one difference. The ground truth RULs for the test set are not revealed. The challenge participants were asked to upload their results (only once per day) to receive a score based on an asymmetrical scoring function (see (Saxena, Goebel, et al., 2008a)). Users can still get their results evaluated using the same scoring function by uploading their results on the repository page, but otherwise it is not possible to compute any other metric on the results in absence of ground truth to allow error computation. The third part of the challenge set is dataset #5V, the final *validation set* that was used to rank the challenge participants, where they were allowed only once chance to submit their results. The challenge since then is still continuing and a participant may submit final results (only once) for evaluation per instructions posted with the dataset on the NASA repository (Saxena & Goebel, 2008).

2.2. Performance Benchmarking

One of the key drivers for this study was to assess state-of-the-art in prognostic methods established through comparisons and performance benchmarking. However, the survey revealed a serious lack of consistency in methods used for performance evaluation. One of the key contributing reasons towards this inconsistency is thought to be the unavailability of established performance benchmark. Originally it was planned that the PHM08 challenge winning performances would establish a benchmark that would allow further improvements as new methods are developed. But since that

Table 1. Description of the five turbofan degradation datasets available from NASA repository.

Datasets	#Fault Modes	#Conditions	#Train Units	#Test Units
Turbofan data from NASA repository	#1	1	100	100
	#2	1	260	259
	#3	2	100	100
	#4	2	249	248
PHM2008 Data Challenge	#5T	1	218	218
	#5V	1	218	435

webpage was taken down in subsequent years these scores have not been easily available except as reported (often partially) in some publications from the winners. It is, therefore, planned to compute several relevant metrics on the submitted results during PHM08 challenge and make them available to serve as reference for future efforts. These benchmarks, however, remain beyond the scope of this paper and will be made available in future publications.

3. C-MAPSS DATASET LITERATURE REVIEW

To analyze various approaches that have been used to solve C-MAPSS dataset problem, all the publications that cite these datasets including the references recommended by the repository were collected through standard web search. The search results returned over seventy publications which were then preprocessed to identify overlapping efforts by same authors or the publications that only cite the dataset but perceptibly did not use them for algorithm development. This resulted in forty unique publications that were then considered for review and analysis in this work.

For the sake of readability, each of these publications were assigned a unique ID to use in various tables summarizing the results presented in this section. This mapping between publication and IDs is presented in Table 10 as appendix. Furthermore, to keep the paper length short, a detailed review analysis of each of the forty publications is not included but only the summarized findings.

The analysis of the collected publications reveals several important observations that are summarized here. First, these publications are binned into various different categories and then analyzed for the distributions thus observed. These categories and corresponding findings are presented next.

3.1. C-MAPSS Dataset Used

Table 2 identifies specific publications that use one or more of these five datasets. It can be observed that the dataset #1 was the most used one (55%), followed by the test set (#5T) from the PHM08 challenge (35%), whereas rest of the other datasets are relatively under utilized. Three publications report generating their own datasets using the C-MAPSS simulator and (Richter, 2012) describes the simulator and how it can be used to generate degradation data rather than using

any specific dataset.

The heavy usage of dataset #1 ($\approx 70\%$) compared to all others among the four from the NASA Repository may be attributed to its simplicity compared to the rest. On the other hand, high usage of dataset #5T is attributed to the PHM08 challenge, where several teams had already used these data extensively, thereby gaining significant familiarity with the dataset as well as developing a natural preference due to availability of corresponding benchmark performance from the challenge leader board.

Table 2. List of publications for each dataset.

Datasets	Publication ID	%	
Turbofan data from NASA repository	#1	5, 6, 10, 13, 14, 15, 19, 20, 23, 24, 25, 26, 27, 28, 31, 32, 33, 34, 36, 37, 38, 40	22/40
	#2	13, 22, 34, 40	4/40
	#3	34, 40	2/40
	#4	7, 34, 40	3/40
PHM08 Data challenge	#5T	1, 2, 3, 4, 8, 12, 16, 17, 21, 29, 30, 34, 35, 40	14/40
	#5V	1, 2, 3, 40	4/40
Simulator	OWN	9, 11, 39	3/40
Other	-	18	1/40

Several publications mentioned in Table 2 have used only the training datasets that have complete (run-to-failure) trajectories. Using data with complete trajectories gives access to the true End-of-Life (EOL) to compute RUL from any time point in a degradation trajectory which could be used to generate a larger set of training data. This approach is also relevant to estimating RULs at different time points and allows the usage of prognostics metrics (Saxena, Celaya, et al., 2008) such as Prognostic Horizon, $\alpha - \lambda$ metric, or the convergence measure. However, in true learning sense the algorithm, once trained, must be tested on unseen data for proper validation, as was required for the PHM'08 challenge datasets. Table 3 shows that 11 different publications used the full training/testing datasets: the training dataset for estimating the parameters of the algorithms and using the full testing datasets for performance evaluation.

3.2. Target Problem Being Solved

As normally expected there is a wide variety of approaches taken in interpreting the datasets, formulating a problem. and

Table 3. List of publications using only *full* training/testing datasets.

Datasets	Publication ID	%	
Turbofan dataset from NASA repository	#1	20, 27, 28, 40	5/40
	#2	40	1/40
	#3	40	1/40
	#4	40	1/40
PHM08 Data challenge	#5T	1, 2, 3, 4, 16, 21, 40	7/40
	#5V	1, 2, 3, 40	4/40

modeling the system to solve the problem. However, contrary to expectations a significant number of publications have utilized these datasets for analysis heavily focused on diagnosis (multi-class classification) rather than prognostics.

By posing a multi-class classification problem various publications attempt to solve mainly three types of problems:

- Supervised classification: The training dataset is labeled (known classes for each feature vector);
- Unsupervised classification: The classes are not known a priori and data are not labeled;
- Partially supervised classification: Some classes are precisely known, others are unknown or are attached with a confidence value to express belief in that class.

Publications 1, 7, 10, 20, 24, 27, 32 use classification for preprocessing steps towards solving a prognostics problem. Specifically, unsupervised classification algorithms are used in publications 1, 7 to segment the dataset into the six operating conditions. For reference, detailed information about various simulated operating conditions in C-MAPSS is described in (Richter, 2012), which can also be used to label these datasets. Supervised and unsupervised classification algorithms are also used in publications 6, 10, 20, 27, 32 to assign a degradation level according to sensor measurements. The sequence of discrete failure degradation stages is indeed relevant for the estimation of the current health state and its prediction (Kim, 2010).

Health assessment, anomaly detection (seen as a 1-class classification problem) or fault identification are tackled in publications 6, 11, 12, 13, 26, 31, 35 using supervised classification methods, and partially supervised classification techniques in publications 12, 27, 33. For these approaches, a known target (or a degradation level) is required to evaluate the classification rate. For instance, four degradation levels were defined for labeling data in publications 6, 10, 27, 33: normal degradation (class 1), knee corresponding to a noticeable degradation (class 2 viewed as a transition between class 1 and 3), accelerated degradation (class 3) and failure (class 4). One such segmentation is provided at URL¹, whereas a different set of segmentation was proposed in publication 13. Using these segmented data (clusters) as proxy to ground

truth, some level of classification performance can be evaluated for comparison purposes.

Similar to several classification approaches used, many approaches were employed for solving the prognostics problem for predicting RULs. In order to give due attention to the analysis of prognostic methods, a discussion is presented separately in Section 4.

3.3. Method for Treatment of Uncertainty

Given the inherent nature of datasets that include several noise factors and lack of specific information on the effects of operational conditions it is important for algorithms to model and account for uncertainty in the system. Different publications have dealt with uncertainty at various stages of processing as described below:

1. **Signal processing step** such as noise filtering using a Kalman filter as in publications 2, 3, 20, Gaussian kernel smoothing in publications 1, 7, and functional principal component analysis in publication 15.
2. **Feature extraction/selection step** such as using principal component analysis and other variants of it as suggested in publications 1, 7, 13, grey-correlation in publication 22, and computing relevance of features for prediction in publication 23.
3. **Health estimation step** such as based on operating conditions assessment to normalize/factor out the effects of operating conditions as proposed in publications 1, 7, 21, 40 and using non-linear regression.
4. **Classification step** where uncertainty modeling plays a role on data labeling using noisy and imprecise degradation levels as shown in publications 12, 27, 33, or on the inference of a sequence of degradation levels such as using Markov Models or multi-models as in publications 6, 10, 24, 32, 34.
5. **Prediction step** such as gradually incorporating prior knowledge during estimation in presence of noise as proposed in publications 4, 14, 16, 17, 19, 21, 30, in determining failure thresholds as in publications 10, 27, 32 or in representing health indicator such as in publication 40 to be used in prediction.
6. **Information fusion step** by merging multiple RUL estimates through Bayesian updating as pointed in publications 4, 21 or in similarity-based matching as in publications 1, 27, 40.

A variety of different uncertainty representation theories are found to be used. Table 4 classifies different publications according to the theory of uncertainty treatment used in corresponding analysis (Klir & Wierman, 1999). As shown in the table, the probability theory is the most popular one (65%) followed by set-membership approaches (in particular fuzzy-sets with 15%), Dempster-Shafer's theory of belief functions

¹<http://members.femto-st.fr/emmanuel-ramasso/data-and-codes>

(13%), and other measures (such as polygon area and Choquet integral).

Table 4. Methods for uncertainty management used on C-MAPSS datasets.

Theories	Publication ID	%
Probability theory	1, 2, 3, 4, 5, 6, 7, 11, 12, 13, 15, 16, 17, 19, 20, 21, 22, 26, 28, 29, 30, 31, 32, 33, 34, 35	26/40
Set-membership	10, 14, 23, 25, 36, 39	6/40
Belief functions	6, 10, 24, 27, 33	5/40
Other measures	10, 40	2/40

3.4. Methods used for Performance Evaluation

Table 5 summarizes the performance measures that have been used for prognostics-oriented publications. A taxonomy of performance measures for RUL estimation was proposed in (Saxena, Celaya, et al., 2008; Saxena, Celaya, Saha, Saha, & Goebel, 2010), where different categories were presented: accuracy-based, precision-based, robustness-based, trajectory-based, computational performance and cost/benefit measures, as well as some measures dedicated specifically to prognostics (PHM metrics). Since this problem involves predictions on multiple units it is expected that the majority of publications would use error-based accuracy and precision metrics. Metric like the Mean Squared Error (MSE) has been used in two different ways- for the estimation of the goodness of fit between a predicted and a real signal, and as an accuracy-based metric to aggregate errors in RUL estimation. Only the publications that fall under latter category are included in the table. The table clearly shows that accuracy-based measures were most widely used, in particular the scoring function from PHM08 challenge, which also weighs accuracy by timeliness of predictions. Broader usage of this metric is also explained by the fact that this is the only metric for which scores from data challenge were available and can be used as benchmark to compare with any new development. However, one may also compute additional measures if using only the training datasets where full trajectories are available. In that case approaches like leave-one-out validation become applicable where all training instances but one are used for training each time and the remaining one is used for performance evaluation. Then the average of the performance measure is computed from all the runs. Publication 27 presents this approach for dataset #1 and a cross-validation procedure for dataset #5T is used in publication 21. Note that publications 19, 20, 32 provide the only RULs estimates for all testing instances (without computing any metrics) and publications 10, 27 present distribution of errors.

4. PROGNOSTIC APPROACHES

C-MAPSS datasets were generated to allow development and benchmarking of various prognostics approaches. However,

Table 5. Performance measures used in prognostics-oriented publications applied on C-MAPSS.

Categories	Measures	Publication ID	%
Accuracy	PHM08 Score	1, 2, 4, 5, 8, 16, 21, 29, 30, 40	10/40
	FPR, FNR	8, 10, 27, 40	4/40
	MSE	3, 8, 15, 17, 29, 40	6/40
	MAPE	4, 23, 28, 32, 34, 39, 40	7/40
	MAE	5, 13, 38, 40	4/40
Precision	ME	25,28,32,39	4/40
	MAD	25	1/40
Prognostics	PH	7, 22	2/40
	$\alpha - \lambda$	7, 22	2/40
	RA	7, 22, 34	3/40
	CV	7, 22, 34	3/40
	AB	34	1/40

as observed from the literature review (see Section 3.2) many researchers have used them to cast a multiclass classification problem instead, even though majority of publications did use them to develop prognostics algorithm. This section focuses on describing those prognostic approaches. These approaches used on C-MAPSS datasets can be divided into three broad categories as described next.

4.1. Category 1: Using functional mappings between set of inputs and RUL

Methods in this category (see Table 6) first transform the training data (trajectories) into a multidimensional feature space and use corresponding RULs to label corresponding feature vectors. Then using supervised learning methods a mapping between feature vectors and RULs is developed. Methods within this category are mostly based on *Neural Networks* with various architectures. Different sensor channels were used to generate corresponding features. However, it was observed that the approaches yielding good performance also included a feature selection step through advanced parameter optimization such as using genetic algorithm and Kalman filtering as described in publications 2, 3 that ranked 2d and 3rd respectively in the competition.

Table 6. Category 1 methods using a mapping learned between a subset of sensor measurements as inputs and RUL as output.

Methods	Publication ID
RNN, EKF	2
MLP, RBF, KF, Ensemble	3
MLP	8
ANN	9
ESN	20
Fuzzy rules, genetic algorithm	36
MLP, adaboost	38

4.2. Category 2: Functional mapping between health index (HI) and RUL

Methods listed in Table 7 are based on the estimation of two mapping functions: One maps sensor measurements to a health index (1-D variable) for each *training unit* based on sensor measurements; The second mapping links health index values to the RUL. These approaches construct a *library of degradation models*. Inference of the RUL for a given test instance includes using the library as *prior knowledge* to *update* the parameters of the model corresponding to the new test instance. Updating can be done using Bayes rule as proposed in publication 4 or other *model averaging* or *ensemble* techniques designed to take into account the uncertainty inherent to the model selection process (Raftery, Gneiting, Balabdaoui, & Polakowski, 2003).

Table 7. Type 2 methods using health index as input and RUL as output.

Methods	Publication ID
Quadratic fit, Bayesian updating	4
Logistic regression	5
Kernel regression, RVM	7
RVM	16
Gamma process	17
Linear, Bayesian updating	19
RVM, SVM, RNN, Exponential and quadratic fit, Bayesian updating	21
Exponential fit	28
Wiener process	29
Copula	30
HMM, LS-SVR	34

Table 8 lists some other approaches that use approximation functions to represent the evolution of individual sensor measurement through time. Given a test instance as many predictions are made as the number of sensors. These predictions are then used in a classifier that assigns a class label related to identified degradation level. Some of these approaches also update classifier parameters with new measurements using some Bayesian updating rules as mentioned previously. These methods were however applied only on dataset #1 in which sensors depict clear monotonic trends.

Table 8. Category 2 methods based on individual sensor modeling and classification.

Methods	Publication ID
exTS, supervised classification	10
SVR	13
exTS, ARX	14
ANN, ANFIS	23
Piece-wise linear (multi-models)	24
exTS	25
ELM, unsupervised classification	32

4.3. Category 3: Similarity-based matching

In these methods (Table 9), historical instances of the system (sensor measurements trajectories labeled with known failure times) are used to create a library. For a given test instance similarity with instances in the library is evaluated generating a set of Remaining Useful Life (RUL) estimates that are eventually aggregated using different methods. Compared to category 2 methods, these methods do not make use of training trajectory abstraction into features, but trajectory data (possibly filtered) are themselves stored. Similarity is computed in the sensor space as in publication 27 or using health indices as in publications 1, 7, 17, 21, 40.

As mentioned in publications 1, 7, in practice the test instance and the training instance may take different time in reaching a particular degradation level from the initial healthy state. Therefore, similarity-based matching must accommodate this difference in the early phases of degradation curves. In publication 40, this problem was tackled by assuming a constant initial wear for all instances yielding an offset on health indices. Efficient similarity measures are also necessary to cope with noise and degradation paths. For instance, in publications 1, 7 three different similarity measures were used, and in publication 40, computational geometry tools were used for instance representation and similarity evaluation.

Table 9. Category 3 methods using similarity-based matching.

Methods	Publication ID
HI-based 3 similarity measures and kernel smoothing	1, 7
Similar to 1 and 7 using 1 similarity measure	22
Feature-based similarity, 1 similarity measure, ensemble, degradation levels classification	27
HI-based similarity, polygon coverage similarity, ensemble	40

An advantage of approaches in this category is that new instances can be easily incorporated. Moreover, similarity-based matching approaches have demonstrated good generalization capability on all C-MAPSS datasets as shown in publications 1, 7, 40 despite a high level of noise, multiple simultaneous fault modes, and a number of operating conditions. This category of algorithms are relatively easily parallelized to reduce computational times needed for inference.

5. SOME GUIDELINES TO USING C-MAPSS DATASETS

Another contribution from this paper is through summarizing some guidelines in using C-MAPSS datasets that may help future users to understand and utilize these datasets better. It summarizes information gathered from the literature review and authors' own experiences, which in many cases goes beyond the documentation provided along with the datasets. Specifically, it offers some general processing steps and lists relevant publications that describe implementation of these

preprocessing steps that could be useful in developing a prognostic algorithm (Figure 1).

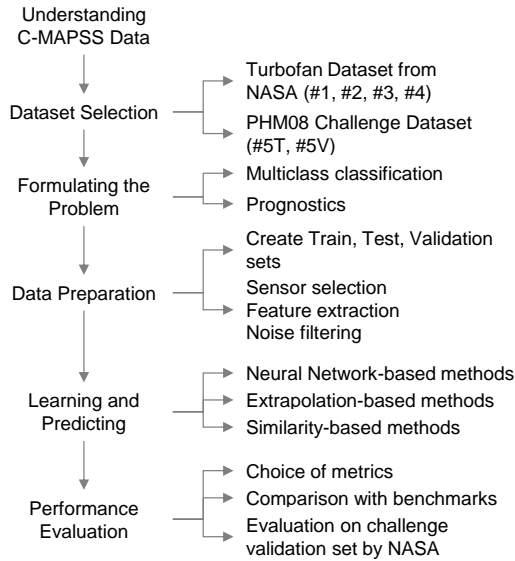


Figure 1. Guidelines to Using C-MAPSS Datasets.

Based on the analysis presented in (Section 3), five general data processing and algorithmic steps are considered:

[Step 1:] Understanding C-MAPSS datasets – Comprehensive background information on turbofan engines and C-MAPSS datasets is well presented in three publications, (Saxena, Goebel, Simon, & Eklund, 2008b), (Richter, 2012), and (T. Wang, 2010). More details about the hierarchical decomposition of the simulated system into critical components can also be found in (Frederick et al., 2007; Abbas, 2010), which provides valuable domain knowledge. These publications do not focus on the physics-of-failure of turbofan engines but describe generation of these datasets and various practical aspects when using C-MAPSS datasets for prognostics. These include description of sensors measurements, illustrations of operating conditions, impact of fault modes, etc., which can play an important role in improving data-driven prognostics algorithms as well. Going from dataset #1 to #4 represents varying degrees of complexity and, therefore, it is recommended to use them in that order to incrementally develop methods to accommodating individual complexity one by one. The challenge datasets fall somewhere in the middle as far as complexity level goes but suffer from availability of ground truth information for a quicker feedback during algorithm development. Therefore, these datasets may be used as validation examples and should be compared to other approaches using benchmarks presented in Section 2.2.

[Step 2:] Defining the problem – Given the nature of these datasets several types of problems can be defined. As men-

tioned in Section 3.2 in addition to prediction, a multi-class classification problem can be defined for a multidimensional feature space. However, the intent behind these data was to promote prognostics algorithm development. Since these data consist of multiple trajectories a problem to predict RUL for all trajectories can be constructed just as the one posed in the data challenge. However, one could also define the problem at a higher granularity by modeling the degradation for each trajectory individually and predict RULs at multiple time instances, which would be more of a condition based prognostics context.

[Step 3:] Data preparation – After a dataset (turbofan or data challenge) is selected, it is suggested to split the original training dataset into two subsets: a training dataset for model parameter estimation (learning) and a testing dataset to test the learned model (see for example publications 21, 40). For the datasets #1 – 4 corresponding RUL vectors are provided for the test sets so users can validate their algorithms. However for the challenge datasets the evaluations can only be obtained by submitting the RULs to the NASA website (on once per day basis for #5T and only once for #5V). Therefore, it may be desirable to split the training set itself for training, test, and validation purposes during algorithm development. The next step is to downselect sensors to reduce problem dimensionality. Some data exploration and preparation approaches for the data challenge (datasets #5T and #5V) are well described in publications 1, 2 and 7. Some “heuristic rules” to avoid over-predictions are also presented in publication 40 and applied on all five C-MAPSS datasets. Some of the better performing methods are based on a PCA such as in publication 1, and other sensor selection procedures such as in publications 2, 3 and 40. From the survey it was noted that the most commonly selected subset of sensors was 7, 8, 9, 12, 16, 17, 20 (as it was also initially suggested in publication 1). Additional sensors may also be considered, similar to the approach proposed in publication 40 where a total of 511 combinations were studied for each dataset for an exhaustive evaluation.

[Step 4:] Learning and Predicting – This step forms the core of prediction problem. As described in Section 3 a variety of learning approaches can be employed to learn various mappings between the sensor data and system health to compute RULs. Some of these methods try to learn RUL as a function of sensor data (system state) or features thereof, others estimate a health index first. Each of the trajectory can be modeled into a degradation process to predict when they cross the zero health threshold using regression methods. Approaches based on health index computation can be applied to all datasets. The approach proposed in publications 1, 7 is the simplest to implement. To deal with normalization (or alternatively segmentation) of data by operating conditions one could use a clustering approach as suggested by the authors above, or one may directly use the parameters described

in publication 18 to validate the performance of segmentation. Some variants for health indicator estimation can also be picked from publications 21 and 40.

[Step 5:] Performance evaluation – Once a learned model results in to satisfactory results on the testing set aside by partitioning the training data, one may use the actual test dataset provided with the datasets. After further tuning, esp. for datasets(#5T and #5V) a final validation can be done by submitting the results to the NASA repository per instructions provided there and receiving the scores. Before uploading the final submission, the generalization capability should be ensured by computing using several performance metrics as discussed in Section 2.2. Some benchmarks have been provided in Section 2.2 using metrics that aggregate prediction performance from multiple units. While the exact numbers would not match, the performance is expected to be in the similar range for results obtained from turbofan datasets that have access to RULs. For comparison purposes, the scores obtained in previous works on complete C-MAPSS trajectories are summarized in publication 40. Note that here using the full trajectory data it is possible to compute prognostics metrics as presented in (Saxena, Celaya, et al., 2008; Saxena et al., 2010) as the actual EOL is known apriori. This allows testing the critical time aspect of a prediction in addition to accuracy and precision measures.

6. CONCLUSION

As observed from published PHM literature the most widely used datasets for data-driven prognostics come from the C-MAPSS turbofan simulator from among the other openly available prognostic datasets. Guided by this observation, a survey of approaches developed using these datasets (since 2008) was carried out with the purpose of understanding the current state-of-the-art and assess how these datasets have helped in development of prognostic algorithms. However, it was noticed that due to several factors these datasets did not get used as intended and any meaningful comparison between approaches was not trivial. Specifically following observations were made and this paper tries to alleviate some of these factors to improve usage of these datasets as originally intended.

- Despite several thousand downloads only 70 papers referring to C-MAPSS were found in the published literature. This suggests that a vast majority of those who downloaded did not get to utilize these data to the point of publishing the results in a publication. Therefore, some guidance has been provided to help in understanding these datasets and how a prognostics problem may be set up in few different ways. Furthermore, a description of all five C-MAPSS datasets is provided identifying their distinguishing characteristics and clearing up some misunderstandings as identified from the survey.
- Among the 70 papers, only a few actually used the testing datasets for evaluating their methods. A mix of different datasets and the metrics used to evaluate performance was observed from the survey. This made it difficult to compare performance between different reported methods in a consistent manner. Therefore, a better explanation of differences in these datasets and providing the top thirty scores from challenge datasets should help future users in comparing their methods against a benchmark in a more consistent manner. Furthermore, it is also suggested how results from datasets that are not from the challenge could be compared against this benchmark established on the challenge set.
- The survey reveals usage of various prognostics approaches that can be divide into three main categories. These approaches are briefly described with potential areas for further improvement. The survey also demonstrated that C-MAPSS datasets can be used for developing and testing methods for several intermediate steps in prognostics such as sensor selection, health indicator estimation, operating conditions modeling in addition to fault estimation and prediction.

With the analysis presented in this paper and references to a variety of approaches employed, this paper hopes to establish public knowledge that can be used by future users in prognostic algorithm development and aid in fulfilling the underlying intent of data repository to facilitate algorithm benchmarking and further development. The issue of performance benchmarking remains to be explored as part of future work where authors plan to compute performance for challenge entries based on several other metrics that will allow comparisons with performance results reported in many publications.

NOMENCLATURE

PHM	Prognostics and Health Management
RUL	Remaining Useful Life
C-MAPSS	Commercial Modular Aero-Propulsion System Simulation
HI	Health index
MLP	MultiLayer Perceptron
ANN	Artificial neural network
RNN	Recurrent neural network
RBF	Radial basis function
ESN	Echo state network
ELM	Extreme learning machine
EKF	Extended Kalman filter
KF	Kalman filter
SVR	Support vector regression
LS-SVR	Least squared support vector regression
exTS	Evolving extended Takagi-Sugeno system

ARX	Autoregressive exogeneous model
ANFIS	Adaptive neuro fuzzy inference system
RVM	Relevance vector machine
HMM	Hidden Markov model
PCA	Principal components analysis
MSE	Mean squared error
MAPE	Mean absolute percentage error
MAE	Mean absolute error
ME	Mean error
PH	Prediction horizon
AP	Acceptable predictions (rate)
$\alpha - \lambda$	Accuracy at specific times
RA	Relative accuracy
CV	Convergence
AB	Average bias
FPR	False positive rate
FNR	False negative rate

ACKNOWLEDGMENT

This work has been partly carried out within the Laboratory of Excellence ACTION through the program “Investments for the future” managed by the French National Agency for Research (ANR-11-LABX-01-01) and with partial support from NASA’s System-wide Safety and Assurance Technologies (SSAT) Project under ARMD/Aviation Safety program.

REFERENCES

- Abbas, M. (2010). *System level health assessment of complex engineered processes*. Unpublished doctoral dissertation, Georgia Institute of Technology.
- Al-Salah, T., Zein-Sabatto, S., & Bodruzzaman, M. (2012). Decision fusion software system for turbine engine fault diagnostics. In *Southeastcon, 2012 proceedings of IEEE* (p. 1-6).
- Coble, J. (2010). *Merging data sources to predict remaining useful life - an automated method to identify prognostic parameters*. Unpublished doctoral dissertation, University of Tennessee, Knoxville.
- Coble, J., & Hines, J. (2008). Prognostic algorithm categorization with phm challenge application. In *Ieee int. conf. on prognostics and health management*.
- Coble, J., & Hines, W. (2011). Applying the general path model to estimation of remaining useful life. *International Journal of Prognostics and Health Management*, 2, 1-13.
- El-Koujok, M., Gouriveau, R., & Zerhouni, N. (2011). Reducing arbitrary choices in model building for prognostics: An approach by applying parsimony principle on an evolving neuro-fuzzy system. *Microelectronics Reliability*, 51(2), 310 - 320.
- Frederick, D., DeCastro, J., & Litt, J. (2007). *User’s guide for the commercial modular aero-propulsion system simulation (C-MAPSS)* (Tech. Rep.). Cleveland, Ohio 44135, USA: National Aeronautics and Space Administration (NASA), Glenn Research Center.
- Gouriveau, R., Ramasso, E., & Zerhouni, N. (2013). Strategies to face imbalanced and unlabelled data in PHM applications. *Chemical Engineering Transactions*, 33, 115-120.
- Gouriveau, R., & Zerhouni, N. (2012). Connexionist-systems-based long term prediction approaches for prognostics. *IEEE Trans. on Reliability*, 61, 909-920.
- Heimes, F. (2008). Recurrent neural networks for remaining useful life estimation. In *Ieee int. conf. on prognostics and health management*.
- Hu, C., Youn, B., Wang, P., & Yoon, J. (2012). Ensemble of data-driven prognostic algorithms for robust prediction of remaining useful life. *Reliability Engineering and System Safety*, 103, 120 - 135.
- Ishibashi, R., & Nascimento Junior, C. (2013). GFRBS-PHM: A genetic fuzzy rule-based system for phm with improved interpretability. In *Ieee conference on prognostics and health management (phm)* (p. 1-7).
- Javed, K., Gouriveau, R., Zemouri, R., & Zerhouni, N. (2012). Features selection procedure for prognostics: An approach based on predictability. In *8th ifac symposium on fault detection, supervision and safety of technical processes* (Vol. 8, p. 25-30).
- Javed, K., Gouriveau, R., & Zerhouni, N. (2013). Novel failure prognostics approach with dynamic thresholds for machine degradation. In *Ieee industrial electronics conference*.
- Jianzhong, S., Hongfu, Z., Haibin, Y., & Pecht, M. (2010). Study of ensemble learning-based fusion prognostics. In *Prognostics and health management conference, 2010. phm '10.* (p. 1-7).
- Kim, H.-E. (2010). *Machine prognostics using health state probability estimation*. Unpublished doctoral dissertation, School of engineering systems, Faculty of built environmental engineering, Queensland university of technology.
- Klir, G., & Wierman, M. (1999). Uncertainty-based information elements of generalized information theory. In (chap. Studies in fuzzyness and soft computing). Physica-Verlag.
- Li, X., Qian, J., & Wang, G. (2013). Fault prognostic based on hybrid method of state judgment and regression. *Advances in Mechanical Engineering*, 2013(149562), 1-10.
- Liao, H., & Sun, J. (2011). Nonparametric and semi-parametric sensor recovery in multichannel condition monitoring systems. *IEEE Transactions on Automation Science and Engineering*, 8(4), 744-753.
- Lin, Y., Chen, M., & Zhou, D. (2013). Online probabilistic operational safety assessment of multi-mode engineering systems using Bayesian methods. *Reliability*

- Engineering & System Safety*, 119(0), 150 - 157.
- Liu, K., Gebraeel, N. Z., & Shi, J. (2013). A data-level fusion model for developing composite health indices for degradation modeling and prognostic analysis. *IEEE Trans. on Automation Science and Engineering*.
- Peel, L. (2008). Data driven prognostics using a Kalman filter ensemble of neural network models. In *Int. conf. on prognostics and health management*.
- Peng, Y., Wang, H., Wang, J., Liu, D., & Peng, X. (2012). A modified echo state network based remaining useful life estimation approach. In *Ieee phm conference*.
- Peng, Y., Xu, Y., Liu, D., & Peng, X. (2012). Sensor selection with grey correlation analysis for remaining useful life evaluation. In *Annual conference of the phm society*.
- Raftery, A., Gneiting, T., Balabdaoui, F., & Polakowski, M. (2003). Using bayesian model averaging to calibrate forecast ensembles. *American Meteorological Society*, 133(5), 1155-1174.
- Ramasso, E. (2009). Contribution of belief functions to hidden Markov models with an application to fault diagnosis. In *Machine learning for signal processing*.
- Ramasso, E. (2014a). Investigating computational geometry for failure prognostics. *Int. Journal on Prognostics and Health Management*. (submitted)
- Ramasso, E. (2014b). Investigating computational geometry for failure prognostics in presence of imprecise health indicator: Results and comparisons on c-mapss datasets. In *European conf. on prognostics and health management*.
- Ramasso, E., & Denoeux, T. (2013). Making use of partial knowledge about hidden states in hidden Markov models: an approach based on belief functions. *IEEE Transactions on Fuzzy Systems*, 10.1109/TFUZZ.2013.2259496.
- Ramasso, E., & Gouriveau, R. (2010). Prognostics in switching systems: Evidential Markovian classification of real-time neuro-fuzzy predictions. In *Ieee prognostics and health management conference*.
- Ramasso, E., & Gouriveau, R. (2013). RUL estimation by classification of predictions: an approach based on a neuro-fuzzy system and theory of belief functions. *IEEE Transactions on Reliability*, *Accepted*.
- Ramasso, E., Rombaut, M., & Zerhouni, N. (2013). Joint prediction of observations and states in time-series based on belief functions. *IEEE Transactions on Systems, Man and Cybernetics - Part B: Cybernetics*, 43, 37-50.
- Riad, A., Elminir, H., & Elattar, H. (2010). Evaluation of neural networks in the subject of prognostics as compared to linear regression model. *International Journal of Engineering & Technology*, 10, 52-58.
- Richter, H. (2012). Engine models and simulation tools. In *Advanced control of turbofan engines* (p. 19-33). Springer New York.
- Sarkar, S., Jin, X., & Ray, A. (2011). Data-driven fault detection in aircraft engines with noisy sensor measurements. *Journal of Engineering for Gas Turbines and Power*, 133, 081602.
- Saxena, A., Celaya, J., Balaban, E., Goebel, K., Saha, B., Saha, S., & Schwabacher, W. (2008). Metrics for evaluating performance of prognostic techniques. In *Int. conf. on prognostics and health management*.
- Saxena, A., Celaya, J., Saha, B., Saha, S., & Goebel, K. (2010). Metrics for offline evaluation of prognostic performance. *International Journal of Prognostics and Health Management*.
- Saxena, A., & Goebel, K. (2008). C-mapss data set. *NASA Ames Prognostics Data Repository*.
- Saxena, A., Goebel, K., Simon, D., & Eklund, N. (2008a). Damage propagation modeling for aircraft engine run-to-failure simulation. In *Prognostics and health management, 2008. phm 2008. international conference on* (pp. 1-9).
- Saxena, A., Goebel, K., Simon, D., & Eklund, N. (2008b). Damage propagation modeling for aircraft engine run-to-failure simulation. In *Int. conf. on prognostics and health management* (p. 1-9). Denver, CO, USA.
- Serir, L., Ramasso, E., & Zerhouni, N. (2012). An evidential evolving multimodeling approach for systems behavior prediction. In *Annual conference of the phm society*.
- Siegel, D. (2009). *Evaluation of health assessment techniques for rotating machinery*. Unpublished master's thesis, Division of Research and Advanced Studies of the University of Cincinnati.
- Son, K. L., Fouladirad, M., & Barros, A. (2012). Remaining useful life estimation on the non-homogenous gamma with noise deterioration based on gibbs filtering : A case study. In *Ieee int. conf. on prognostics and health management*.
- Son, K. L., Fouladirad, M., Barros, A., Levrat, E., & Iung, B. (2013). Remaining useful life estimation based on stochastic deterioration models: A comparative study. *Reliability Engineering and System Safety*, 112, 165 - 175.
- Sun, J., Zuo, H., Wang, W., & Pecht, M. (2012). Application of a state space modeling technique to system prognostics based on a health index for condition-based maintenance. *Mechanical Systems and Signal Processing*, 28, 585 - 596.
- Tamilselvan, P., & Wang, P. (2013). Failure diagnosis using deep belief learning based health state classification. *Reliability Engineering and System Safety*, 115(0), 124 - 135.
- Wang, P., Youn, B., & Hu, C. (2012). A generic probabilistic framework for structural health prognostics and uncertainty management. *Mechanical Systems and Signal Processing*, 28, 622 - 637.
- Wang, T. (2010). *Trajectory similarity based prediction for remaining useful life estimation*. Unpublished doctoral

dissertation, University of Cincinnati.

- Wang, T., Yu, J., Siegel, D., & Lee, J. (2008). A similarity-based prognostics approach for remaining useful life estimation of engineered systems. In *Int. conf. on prognostics and health management* (p. 1-6).
- Xi, Z., Jing, R., Wang, P., & Hu, C. (2013). A copula-based sampling method for data-driven prognostics and health management. In *Asme 2013 international design engineering technical conferences and computers and information in engineering conference*.
- Xue, Y., Williams, D., & Qiu, H. (2011). Classification with imperfect labels for fault prediction. In *Proceedings of the first international workshop on data mining for service and maintenance* (pp. 12–16). ACM.
- Yu, J. (2013). A nonlinear probabilistic method and contribution analysis for machine condition monitoring. *Mechanical Systems and Signal Processing*, 37, 293-314.
- Zein-Sabatto, S., Bodruzzaman, J., & Mikhail, M. (2013). Statistical approach to online prognostics of turbine engine components. In *Southeastcon, 2013 proceedings of IEEE* (p. 1-6).
- Zhao, D., P., R. G., & Willett. (2011). Comparison of data reduction techniques based on SVM classifier and SVR performance. In *Proc. SPIE, signal and data processing of small targets* (Vol. 8137, p. 1-15).

APPENDIX

All references were mapped to numeric identifiers to be used in survey and analysis results for better readability. This mapping is provided in the Table 10 below.

Table 10. References to ID mapping.

Reference	Publication ID
(T. Wang, Yu, Siegel, & Lee, 2008)	1
(Heimes, 2008)	2
(Peel, 2008)	3
(Coble & Hines, 2008)	4
(Coble, 2010)	
(Coble & Hines, 2011)	5
(Siegel, 2009)	
(Ramasso, 2009)	6
(T. Wang, 2010)	7
(Riad, Elminir, & Elattar, 2010)	8
(Abbas, 2010)	9
(Ramasso & Gouriveau, 2010)	10
(Ramasso & Gouriveau, 2013)	
(Sarkar, Jin, & Ray, 2011)	11
(Xue, Williams, & Qiu, 2011)	12
(Zhao, P., & Willett, 2011)	13
(El-Koujok, Gouriveau, & Zerhouni, 2011)	14
(Liao & Sun, 2011)	15
(P. Wang, Youn, & Hu, 2012)	16
(Son, Fouladirad, & Barros, 2012)	17
(Richter, 2012)	18
(Sun, Zuo, Wang, & Pecht, 2012)	19
(Peng, Wang, Wang, Liu, & Peng, 2012)	20
(Hu, Youn, Wang, & Yoon, 2012)	21
(Peng, Xu, Liu, & Peng, 2012)	22
(Javed, Gouriveau, Zemouri, & Zerhouni, 2012)	23
(Serir, Ramasso, & Zerhouni, 2012)	24
(Gouriveau & Zerhouni, 2012)	25
(Yu, 2013)	26
(Ramasso, Rombaut, & Zerhouni, 2013)	27
(Liu, Gebraeel, & Shi, 2013)	28
(Son, Fouladirad, Barros, Levrat, & Lung, 2013)	29
(Xi, Jing, Wang, & Hu, 2013)	30
(Lin, Chen, & Zhou, 2013)	31
(Javed, Gouriveau, & Zerhouni, 2013)	32
(Ramasso & Denoeux, 2013)	33
(Li, Qian, & Wang, 2013)	34
(Tamilselvan & Wang, 2013)	35
(Ishibashi & Nascimento Junior, 2013)	36
(Gouriveau, Ramasso, & Zerhouni, 2013)	37
(Jianzhong, Hongfu, Haibin, & Pecht, 2010)	38
(Zein-Sabatto, Bodruzzaman, & Mikhail, 2013)	39
(Al-Salah, Zein-Sabatto, & Bodruzzaman, 2012)	
(Ramasso, 2014b)	40
(Ramasso, 2014a)	