

Data-Driven Open-Set Fault Classification of Residual Data Using Bayesian Filtering

Daniel Jung

The self-archived postprint version of this journal article is available at Linköping University Institutional Repository (DiVA):

<http://urn.kb.se/resolve?urn=urn:nbn:se:liu:diva-168862>

N.B.: When citing this work, cite the original publication.

Jung, D., (2020), Data-Driven Open-Set Fault Classification of Residual Data Using Bayesian Filtering, *IEEE Transactions on Control Systems Technology*, 28(5), 2045-2052.

<https://doi.org/10.1109/TCST.2020.2997648>

Original publication available at:

<https://doi.org/10.1109/TCST.2020.2997648>

Copyright: Institute of Electrical and Electronics Engineers

<http://www.ieee.org/index.html>

©2020 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.



Data-driven Open Set Fault Classification of Residual Data using Bayesian Filtering

Daniel Jung

Abstract—Data-driven fault classification in industrial applications is complicated by unknown fault classes and limited training data. In addition, different faults can have similar effects on sensor outputs resulting in fault classification ambiguities, i.e. multiple fault hypotheses can explain the data. One solution is to identify and rank all plausible fault classes which give useful information, for example at a workshop when performing troubleshooting. A probabilistic fault classification algorithm is proposed for residual data classification combining Weibull-calibrated one-class support vector machines for fault class modeling and Bayesian filtering for time-series analysis. The fault classifier ranks different fault classes and can identify sequences from unknown fault realizations, i.e. faults not represented in training data. Real residual data computed from sensor data and model analysis of an internal combustion engine are used as a case study illustrating the usefulness of the proposed method.

Index Terms—Fault classification; Open set classification; Machine Learning; Support Vector Machines; Hybrid fault diagnosis.

I. INTRODUCTION

IN the automotive industry, on-board diagnosis (OBD) systems have been used for emission-related monitoring for decades. New applications, such as predictive maintenance and assisted troubleshooting at the workshop, are important to improve reliability and reduce system down-time in order to increase customer value. Connected vehicles and cloud computation capacities have put focus on machine learning methods for fault diagnosis and prognostics.

Fault diagnosis of industrial systems is often conducted by analysis and classification of time-series data collected during system operation, for example sensor data or computed residuals [1], [2]. When designing a fault diagnosis system, there are often many different types of faults that can occur in the system and should be detected. Even though there are tools to systematically identify all these fault classes early in the system development phase, see for example [3], it is still a difficult task, especially for large-scale or complex systems. Therefore, there can be unknown faults that are not taken into consideration when training the diagnosis system [4].

Another complicating factor in data-driven fault diagnosis is collecting representative training data from all relevant faults. Data collection is an expensive and time-consuming process and not feasible in many applications [5], [6]. Especially, since many faults do not occur until after years of operation. Therefore, training data are not representative of all fault scenarios which means that a diagnosis system must be able to identify both known and unknown fault scenarios.

Different faults can have similar effects on system dynamics resulting in fault classification ambiguities. Therefore, it is not desirable that a data-driven classifier only selects one fault class, since the true fault could be missed, but should instead identify and rank all plausible fault classes [4]. This type of information is useful, for example, at the workshop to support a technician during troubleshooting [7]. For reliable fault classification, it is also necessary to identify data sets with unknown faults, i.e., fault scenarios not represented in training data, since these cases need special attention to improve classification performance over time [8].

A. Problem Formulation

The objective of this work is to develop a data-driven fault classification algorithm for time-series data, for example sensor data or model-based residuals, that identifies and ranks fault hypotheses (fault classes). It is assumed that training data are limited and not representative of all fault realizations. Machine learning algorithms that assume all data classes are known and representative training data are available, are not expected to give reliable outputs, especially in fault scenarios where data deviate too much from the training data. A fault classifier should therefore be able to identify when there are data sequences with unknown fault scenarios, i.e. sequences that do not resemble training data.

Fault diagnosis of an internal combustion engine is used as a case study. The fault scenarios cover different types of engine faults, including sensor faults, leakages, and air filter clogging. As input to the data-driven fault classification algorithm, a set of residual data is computed from a physically based model and real data from different fault scenarios collected from the engine test rig [9], see Fig. 1.

B. Related Research

Data-driven monitoring and fault diagnosis of internal combustion engines is investigated in, for example, [10]. A data-driven classifier approach for fault diagnosis of an electric throttle control system is proposed in [5] where incremental learning is used to improve classification performance over time. In [11], an ensemble approach for automotive fault classification of both known and unknown faults in time-series data is developed by combining multiple machine learning methods for classification. A two step fault classification approach to handle unknown faults in an electronic system using Gaussian mixture models and k-means is proposed in [12]. With respect to the mentioned work, an incremental

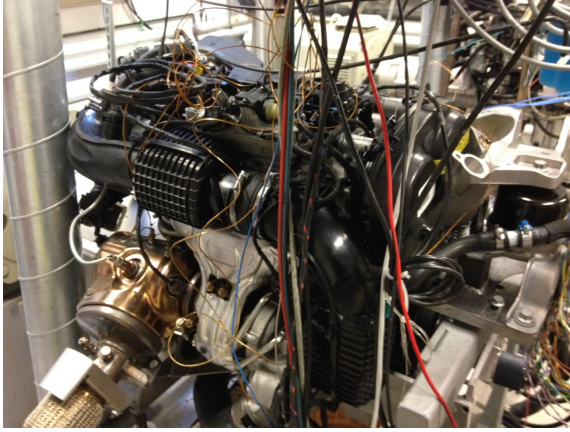


Fig. 1. The picture shows the engine test bench that is used for data collection.

probabilistic fault classification method is proposed that ranks different faults using model-based residuals as input.

One solution to limited training data is to use a physically based model of the system to generate features, for example residuals [4]. Fault diagnosis methods for automotive applications, combining model-based and data-driven methods, are proposed in, for example, [13], [14], [15], [16]. Research highlights the benefits of bridging and combining model-based and data-driven methods for fault diagnosis instead of only focusing on one of them [17].

In [18], both sensor data and residual data are used as input to a tree augmented naive Bayes fault classifier. In [6], a conditional Gaussian network is proposed to handle both known and unknown fault classes. In [19], feature selection using neural networks is applied before training the fault classifiers. In [4], a hybrid diagnosis system design is proposed which combines model-based fault isolation with support vector data description anomaly classifiers to rank the different fault hypotheses. In [20], model-based residuals and sensor data are used as inputs to a Bayesian network to perform fault classification and in [21] model data features are extracted and fed into a neural network classifier. In [22], a hybrid approach combining model-based residuals with hidden Markov models and Bayesian methods is used to classify unknown faults.

Another related research topic is the open set recognition problem in computer vision where data can belong to unknown classes not covered by the training data [8]. Unknown classes are further categorized into *known unknowns* and *unknown unknowns*, where the second case corresponds to the unknown faults considered in this work. Different algorithms have been proposed to solve the open set recognition problem, for example Weibull-calibrated support vector machines [23] and extreme value machines [24].

This work is based on previous research in [4], [25]. The main contribution, with respect to mentioned works, is a data-driven probabilistic classification algorithm of time-series data combining Weibull-calibrated one-class Support Vector Machines [26] and Bayesian filtering and smoothing [27] to improve classification performance and ranking of fault hypotheses.

II. FAULT CLASS MODELING USING OPEN SET CLASSIFICATION

In real-life applications where training data are limited, it is important that a classifier can identify residual data that cannot be explained by any of the known fault classes, i.e. data that significantly deviates from training data.

A. Using One-class Classifiers for Modeling Fault Classes

Let m be the number of available residuals and $\bar{r} = (r_1, r_2, \dots, r_m)$ is a sample of all residuals outputs. The purpose of modeling different fault classes is to identify which fault hypotheses can explain the observed data \bar{r} . One-class classifiers are suitable for modeling fault classes since each class can be modeled individually.

There are multiple methods proposed for one-class classification, for example probabilistic models, one-class support vector machines (OSVM), and isolation forests (iForests) [28]. Probabilistic models use probability distributions to model data from one class and detect outliers, with respect to that class, when the likelihood of a sample is small, see for example [6]. Non-probabilistic models, such as OSVM and iForests, model a decision boundary that encapsulates training data to determine if new data can be explained by that class or not.

Since training data are assumed to be limited, the distribution of data is not expected to be representative of each fault class. Training data might have been collected through experiments to cover different fault realizations but not to be representative of the actual distribution of fault realizations. The objective is to identify plausible fault hypotheses, regardless of how likely they are. Therefore, a non-probabilistic approach is used to model which observations \bar{r} can be explained by each fault class. Training data from each known class is modeled using a decision function representing the maximum distance from any training data where a new sample could be explained by that class, called a *compact abating probability* (CAP) model [23]. Unknown fault classes are identified when data are significantly deviating from training data. Fig. 2 illustrates a set of CAP models and the problem of classifying a set of new data when it significantly deviates from the known fault classes. It is shown in [23] that an OSVM classifier with a radial basis function (RBF) kernel yields a CAP model.

B. One-class Support Vector Machines

There are two similar approaches for designing OSVM classifiers, referred to as ν -SVM [29] and Support Vector Data Description (SVDD) [30], respectively, where ν -SVM is used in this work. An OSVM classifier uses the kernel trick to model a decision boundary that encapsulates data from that class [31]. This is illustrated in Fig. 2 where the black lines represent the decision boundaries of two OSVM classifiers modeling data from Class 1 and Class 2, respectively.

The OSVM classifier computes a score function, when evaluating each new sample, that is positive when belonging to the nominal class or negative if it is considered an outlier, i.e. not belonging to that class. The OSVM classifier evaluates

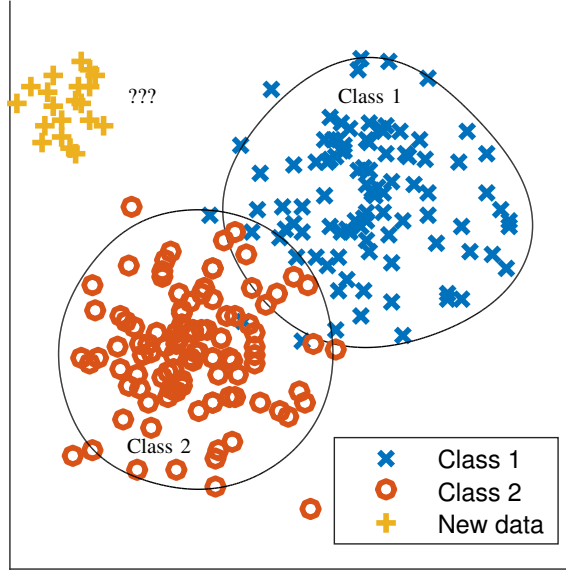


Fig. 2. An illustration of two CAP models using OSVM classifiers to model two known fault classes. The new data cannot be explained by any of the known fault classes and is considered to belong to an unknown fault class.

each sample of residual data independently, meaning that time-series information of the residuals are ignored.

In previous work [4], a set of OSVM classifiers is used to model the CAP models from known fault classes. When classifying new data, each fault class is ranked based on how many samples that are associated to that fault class. Note that a sample can be explained by multiple fault classes. As new data are collected from different faults and correctly classified, the OSVM classifiers are updated accordingly to improve performance over time. Incremental training can be applied to reduce the computational cost when new data are collected, see for example [32].

C. Weibull-calibrated OSVM

Even though the OSVM classifier is a CAP model, its decision boundary depends on the distribution of the support vectors. Therefore, it is relevant to have a measure of the probability that a sample \bar{r} , can be explained by fault class f^l , here denoted $P(\bar{r} \in f^l)$. There are some proposed methods to translate the score computed by a SVM into a probability, for example Platt scaling [33] or Weibull-calibrated SVM [23]. The advantages of Weibull-calibrated SVM with respect to Platt scaling are discussed in, e.g. [23]. However, the Weibull-calibrated SVM classifier is a multi-class classifier that outputs one fault class, which could be the unknown class. Since the objective here is to model each fault class separately, to identify all plausible fault hypotheses, a Weibull-calibrated OSVM method, proposed in [26], is used called P_I -OSVM.

In [23], [26], statistical extreme value theory is applied when proposing the P_I -OSVM classifier to model data from each class. The output scores from the support vectors of the OSVM are modeled to be reverse Weibull distributed. The corresponding cdf of the reverse Weibull distribution measures the probability that a new sample can be explained by that

fault class, referred to as *probability of inclusion* in [26]. An example is shown in Fig. 3 showing the distribution of the OSVM score for a set of data, a reverse Weibull distribution fitted to the score values, and the corresponding cdf.

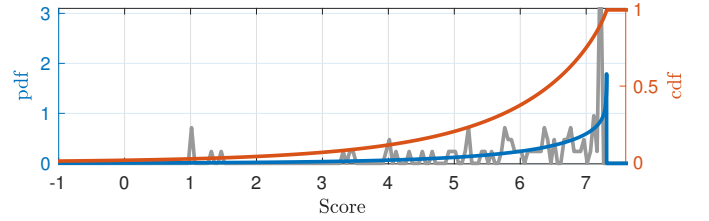


Fig. 3. PDF and CDF of the parameterized reverse Weibull distribution fit to the score values of the support vectors of a OSVM classifier.

The reverse Weibull cdf parameterized for the OSVM score value $g(\bar{r})$ is given by

$$P(\bar{r} \in f^l) = \begin{cases} e^{-\left(\frac{-g(\bar{r}) + \nu_l}{\lambda_l}\right)^{\kappa_l}} & \text{if } -g(\bar{r}) + \nu_l \geq 0 \\ 1 & \text{otherwise} \end{cases} \quad (1)$$

where $\nu_l, \lambda_l, \kappa_l \geq 0$ are fitted parameters for fault class f^l . For (1), denoted P_I -OSVM, to be a CAP model, the probability $P_{f^l}(\bar{r})$ is thresholded by a parameter δ which represents when the Euclidean distance from a new sample to the training data is too large. An example of P_I -OSVM models $P(\bar{r} \in f^l)$ for a set of fault modes f^l are shown for a two residual output case in Fig. 4 where different fault scenarios in training data result in different residual outputs. The z-axis represents the conditional probability (1) that each fault class can explain the residual outputs.

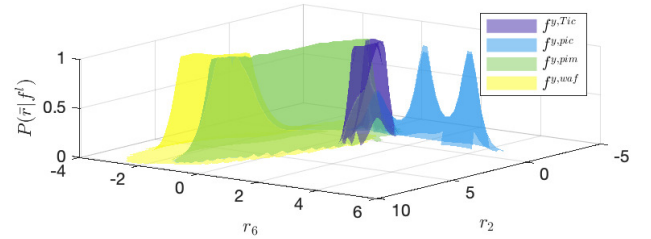


Fig. 4. A set of P_I -OSVM models are parameterized for two residuals. The figure shows the probability of inclusion for each fault class.

III. FAULT CLASSIFICATION OF TIME-SERIES DATA USING BAYESIAN FILTERING AND SMOOTHING

One approach to classify each sample \bar{r}_t , where subscript t is used to denote time index, using the set of P_I -OSVM models is to select the class f^l with the highest probability, i.e., $\arg \max_{f^l} P(\bar{r} \in f^l)$. The probability of a sample to belong to an unknown fault class, denoted f^x , is difficult to model without prior information. In [6], [23], there are no probability models of the unknown fault class f^x . Instead, f^x is selected when the probabilities of all known fault classes are below some threshold. Here, $P(\bar{r} \in f^x) = \delta$ is modeled equal to the threshold of the corresponding CAP models, i.e. samples that do not belong to a known fault class are more likely to come from an unknown fault class.

Since a fault is often present during a longer time interval, Bayesian filtering and smoothing are applied here to improve classification performance by weighing in information from consecutive samples [27]. This is relevant if there are multiple fault classes that can explain the same observations.

The probability that the system is changing from one fault mode to another at time t is modeled using a transition matrix $\Pi \in \mathbb{R}^{n+1 \times n+1}$, where n is the number of known fault classes and plus one for the unknown fault class. Let $\Pi_{l,k}$ denote the element representing the probability that the system changes from mode f_{t-1}^l to f_t^k at time t . Faults are rare events and the probability that the system is changing mode is considered small compared to the system staying in the same mode.

The pdf $p(\bar{r}_t|f^l)$ of the residual output \bar{r}_t given fault class f^l is unknown. However, to be able to use the P_I -OSVM models in a Bayesian framework, it is assumed here that $p(\bar{r}_t|f^l)$ is large when $P(\bar{r} \in f^l)$ is large. Then, the pdf is modeled as $p(\bar{r}_t|f^l) \propto P(\bar{r} \in f^l)$.

A Bayesian filter evaluating the probability of each fault class f_t^l at time t can be computed sequentially as

$$p(f_t^l|\bar{r}_{1:t}) \propto p(\bar{r}_t|f_t^l) \sum_{k=1}^{n+1} \Pi_{k,l} p(f_{t-1}^k|\bar{r}_{1:t-1}) \quad (2)$$

where the prior distribution $p(f_0^l|\bar{r}_0) = p(f_0^l)$ and the probabilities of all modes are normalized, i.e. $\sum_{k=1}^{n+1} p(f_t^k|\bar{r}_{1:t}) = 1$.

The sequential formulation of Bayesian filtering is suitable for on-line computations where class probabilities are computed based on previous samples. A workshop would be able to download logged data and perform off-line computations on the whole data batch. Bayesian smoothing can be applied to a batch of T samples by performing an additional backward filtering after (2) as

$$p(f_t^l|\bar{r}_{1:T}) \propto p(f_t^l|\bar{r}_{1:t}) \sum_{k=1}^{n+1} \Pi_{l,k} p(f_{t+1}^k|\bar{r}_{1:T}) \quad (3)$$

followed by a normalization to compute the class probabilities. Combining the P_I -OSVM classifiers and Bayesian filtering or smoothing gives a systematic method to identify and rank the different fault hypotheses based on how many samples in a data batch are classified as each fault class [4].

IV. CASE STUDY

The case study in this work is the same internal combustion engine system as considered in [25] and [9]. Sensor data have been collected from the engine test bed, including nominal system behavior (NF - No Fault) and seven different single-fault scenarios: air filter clogging f^{paf} , leakages at the air filter f^{Waf} and at the throttle f^{Wth} , and four different sensor faults $f^{y,Tic}$, $f^{y,pic}$, $f^{y,pim}$, and $f^{y,Waf}$. Table I summarizes the seven fault scenarios. The locations of the four sensors are shown in Fig. 5 where y^{Tic} and y^{pic} measure the temperature and pressure after the intercooler, y^{pim} measures the pressure at the intake manifold, and y^{Waf} measures the air flow through the air filter.

A mathematical model is available describing the air flow through an internal combustion engine. The model has been

TABLE I
A SUMMARY OF FAULT SCENARIOS COLLECTED FROM ENGINE TEST RIG.

Fault	Description
f^{paf}	Air filter clogging
f^{Waf}	Leakage after air filter
f^{Wth}	Leakage before throttle
$f^{y,Tic}$	Intermittent fault in sensor measuring temperature at intercooler
$f^{y,pic}$	Intermittent fault in sensor measuring pressure at intercooler
$f^{y,pim}$	Intermittent fault in sensor measuring intake manifold pressure
$f^{y,Waf}$	Intermittent fault in sensor measuring air flow through air filter

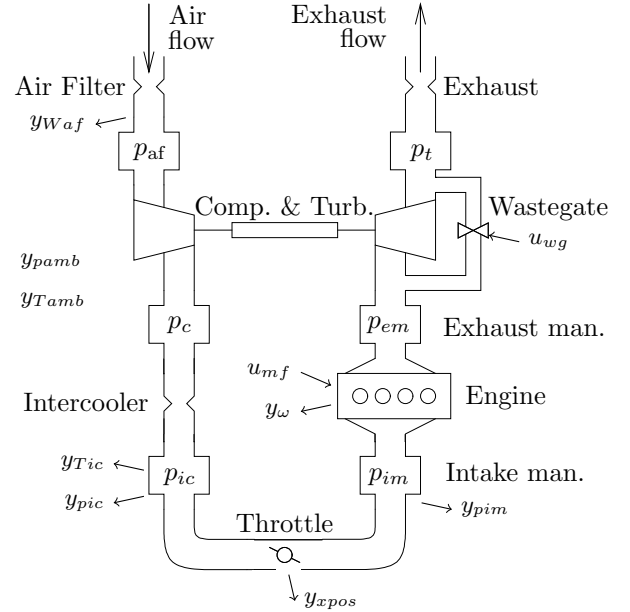


Fig. 5. A schematic of the model of the air flow through the model. This figure is used with permission from [34].

used in previous works for residual generation, see for example [4], and the model structure is similar to the model described in [35], which is based on six control volumes and mass and energy flows given by restrictions, see Fig. 5.

Nine residual generators $\bar{r} = (r_1, \dots, r_9)$ have been generated in [25] from the model, using the Fault Diagnosis Toolbox in Matlab [36]. A residual is a function comparing two different estimates of the same quantity to detect inconsistencies, for example, between a sensor value and a model prediction of the measured quantity. An illustrative example is shown in Fig. 6 where u represents control signals, f faults, y sensor data, \hat{y} model predictions, and $r = y - \hat{y}$ is the residual.

The internal combustion engine is an example of a system that operates at many different operating conditions including transients. The residuals are designed to, ideally, filter out the system dynamics while being sensitive to faults. Even though both sensor data and residuals can be used as inputs to a classifier, only residual data will be used here.

The nine residuals are evaluated using data from different fault scenarios collected from the engine test rig¹. The data set

¹Residual data are available in the Fault Diagnosis Toolbox [36] that can be downloaded from <https://faultdiagnostictoolbox.github.io>. The selected residual subset used in this work is described in [25].

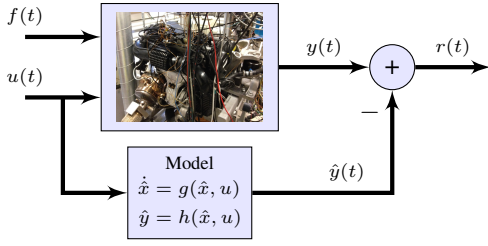


Fig. 6. An example of a residual $r(t)$ comparing measurements from the system $y(t)$ with model predictions $\hat{y}(t)$.

contains 20 276 samples including nominal and faulty data. To evaluate the situation with limited training data, only 10% of the residual data, both nominal operation and from different fault scenarios, are used as training data and the remaining set is used for validation. Figure 7 shows data from each residual from both nominal (NF) and seven fault classes (blue data) and the corresponding fault label (red data). The air filter clogging f^{paf} and leakages f^{Waf} and f^{Wth} have been collected from persistent fault scenarios, while sensor fault data are collected from intermittent sensor faults as shown in the figure.

V. EXPERIMENTAL RESULTS

A P_I -OSVM model (1) is calibrated for each class in the training data and a decision threshold $\delta = 5\%$ is selected for each model. The OSVM classifier, used in the P_I -OSVM models, is implemented using the function `fitcsvm` in Matlab and its kernel parameters are fit to training data using a subsampling heuristic [37]. In this analysis, fault detection and classification are performed simultaneously and the fault-free class NF is included as a fault class.

Validation data from each fault scenario in Fig. 7 are used to evaluate the similarity between the models by analyzing how many samples can be explained by each fault class. Figure 8 shows the percentage of data from each fault scenario that can be explained by each fault class. Samples that are not associated to any known fault class are classified as the unknown fault class f^x . Note that the sum of each column in Fig. 8 can exceed 100% since each sample can belong to multiple fault classes. A significant number of samples can be explained by more than one fault class, e.g. $\{NF, f^{paf}\}$ and $\{f^{waf}, f^{wth}\}$, showing that the CAP models for the different fault classes are overlapping. It is also visible that the overlap is not symmetric between fault pairs. For example, 81% of the samples from fault scenario f^{paf} can also be explained by $f^{y.pim}$ but only 31% of the samples from $f^{y.pim}$ can be explained by f^{paf} .

The CAP models are useful to identify fault hypotheses, i.e. which fault classes could explain the residual data. However, each sample is classified independently of the others ignoring information from the time-series data. To improve fault classification performance, the next step is to take time-series information into consideration.

A. Classification Using Bayesian Filtering and Smoothing

The next step is to evaluate the benefits of applying Bayesian filtering and smoothing with respect to only sample-

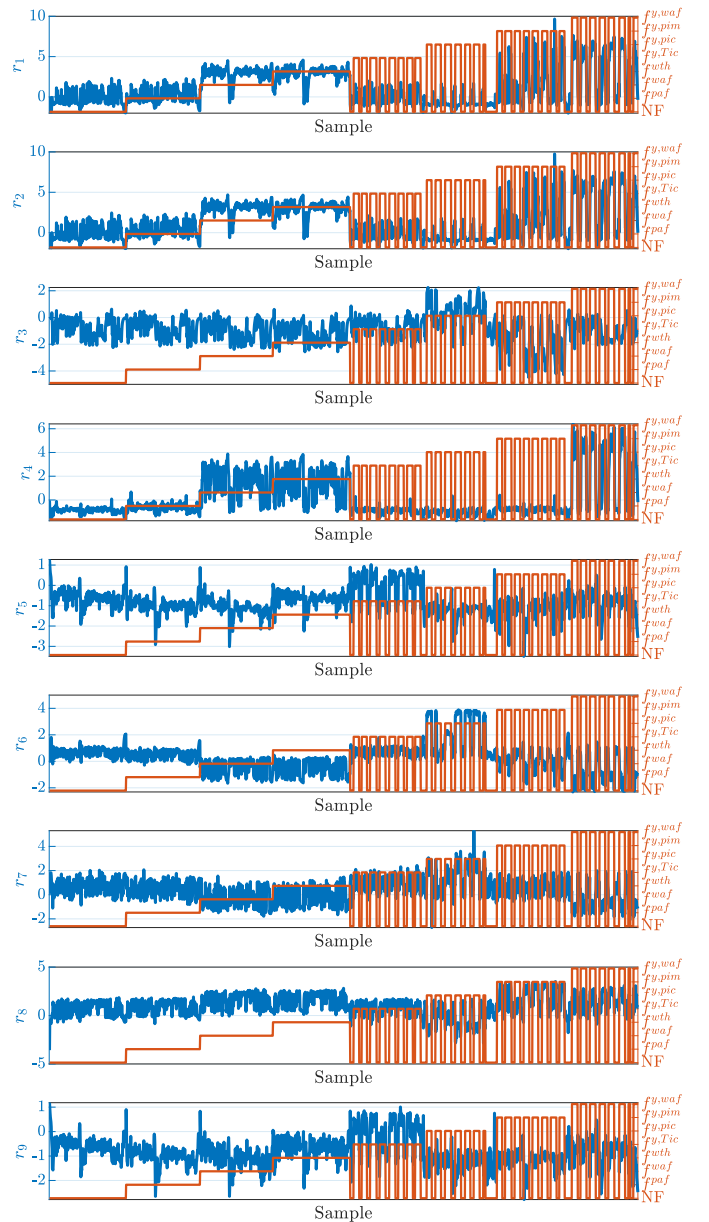


Fig. 7. Data from nine residuals collected from nominal system behavior (NF) and seven different faults. Each subplot shows one residual output where the blue curve is the residual output and the red curve is class label.

by-sample classification of residual data. First, sample-by-sample classification is performed where each sample \bar{r}_t is associated to the fault class f^l with the highest probability $p(f^l|\bar{r}_t)$ at time t . Each fault class f^l is ranked during a fault scenario by counting how many samples are associated with that fault class, similar to what is used in [4] and [25].

The distribution $p(f^l|\bar{r}_t)$ is evaluated using validation data where the a priori distribution $p(f_0^l)$ of all fault classes f^l are assumed equal and the results are shown in Fig. 9. It is highlighted in gray when each fault class is the true fault in the data set. Ideally, $p(f^l|\bar{r}_t) = 1$ when f^l is the true fault class and zero otherwise. Fig. 10 summarizes classification performance when each sample \bar{r} is classified as the fault class f^l with the highest probability $p(f^l|\bar{r}_t)$. Each element (k, l)

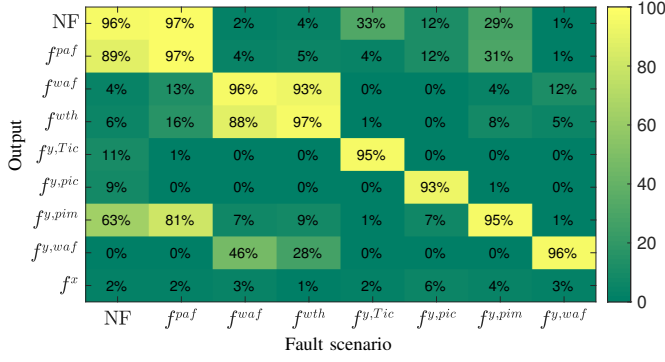


Fig. 8. Modeling fault classes using residual data from Fig. 7 and CAP models, here thresholded P_I -OSVM models. The overlap between fault classes is evaluated by counting the percentage of data that can be explained by each fault class. Samples not belonging to any known fault class belong to the unknown fault class f^x .

in the matrix shows how many samples from a fault scenario with fault f^l are associated to fault class f^k . The evaluation in Fig. 8 shows that it is more difficult to correctly classify fault classes when the CAP-models are overlapping, in this case mainly $\{NF, f^{paf}\}$ and $\{f^{waf}, f^{wth}\}$, respectively.

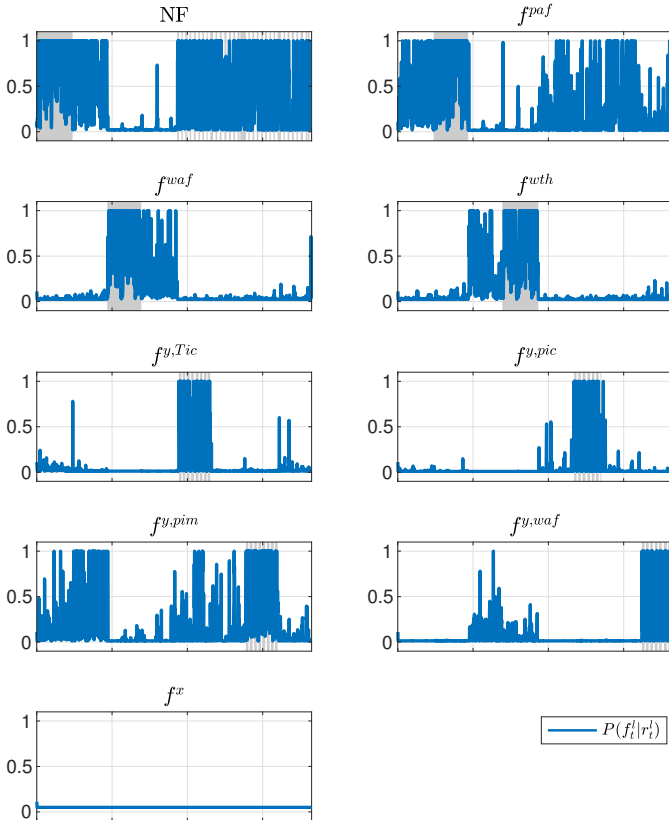


Fig. 9. Fault class probabilities $p(f_t^l | \bar{r}_t)$ from validation data in Fig. 7. The gray intervals represent when the corresponding fault class is the true one and the probability should be high, and zero otherwise.

A comparison of filtered (2) and smoothed (3) estimates of class probabilities are shown in Fig. 11. Here, the transition probability between two different classes in Π is chosen experimentally as 1% and staying at the same class as

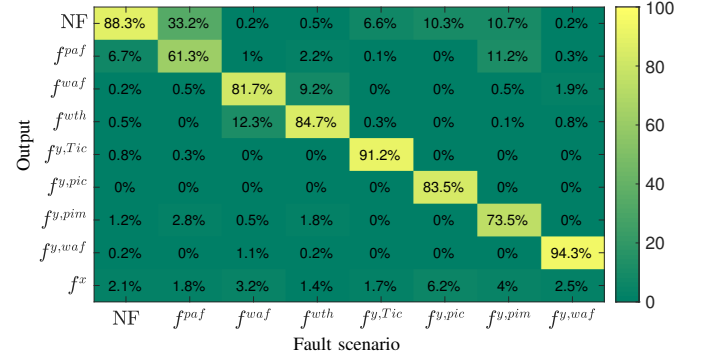


Fig. 10. Evaluation of a set of P_I -OSVM classifiers where the output of the ensemble classifier for each sample is the class with the highest probability, see Fig. 9. Each column represent a fault scenario and each row the ranking of each fault class.

$100 - (n + 1)\%$. Experiments show that a higher transition probability between fault classes results in bigger fluctuations in $p(f_t^l | \bar{r}_t)$ while a lower transition probability reduces the fluctuations but, sometimes, also requires more samples after a fault occur before $p(f_t^l | \bar{r}_t)$ changes significantly. The different subplots in Fig. 11 show the computed probability of each fault class where the highlighted gray areas show when the fault is present and the ranking should be high, and zero otherwise.

Each sample is associated to the fault class with the highest probability after applying Bayesian filtering and smoothing. Compared to the sample-by sample classification in Fig. 9 the filtered estimates significantly improve fault classification performance, especially between fault classes that are overlapping in Fig. 8. The smoothed probability often seems to dominate for one class at each sample time compared to using the Bayesian filter only. In the figure, only a few samples are classified to belong to the unknown fault case.

Classification performance using Bayesian filtering and smoothing are shown in Fig. 12 and Fig. 13, respectively. The output percentages show the ranking of each fault class in each scenario. The most significant improvement, with respect to the sample-to-sample classification in Fig. 8, is classification of fault f^{paf} where the ranking of the true fault increases from 61.3% to 81.3%. When comparing the results in Fig. 12 and Fig. 13, Bayesian smoothing has only a slight improvement in classification accuracy with respect to Bayesian filtering.

B. Classification of unknown faults

Unknown fault scenarios are simulated by training a set of P_I -OSVM models without including training data from the fault class that is considered unknown in the scenario. Seven unknown fault scenarios are evaluated where data from one fault class in Table I are excluded during each training phase and a set of P_I -OSVM models is trained based on the remaining known fault classes. Then, validation data from the unknown fault class is classified using the P_I -OSVM models and Bayesian smoothing to rank the different fault classes in each scenario. Ideally in each fault scenario, the unknown fault class f^x should have the highest rank since the model of the true fault is not available.

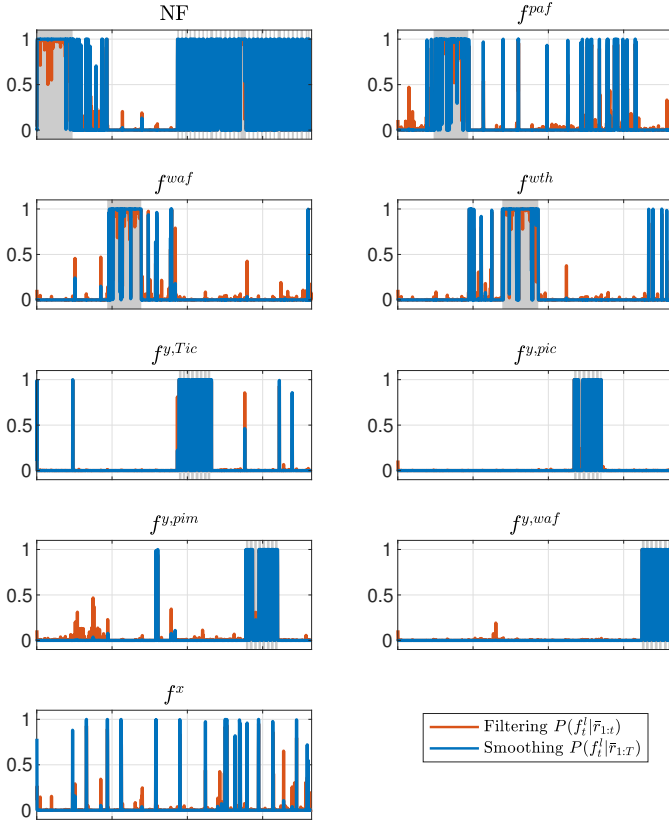


Fig. 11. Fault class probabilities using Bayesian filtering $p(f_t^l | \bar{r}_{1:t})$ and smoothing $p(f_t^l | \bar{r}_{1:T})$. The gray intervals represent when the corresponding fault class is present. Smoothing makes the probability of one fault class more dominating with respect to the other classes compared to filtering.

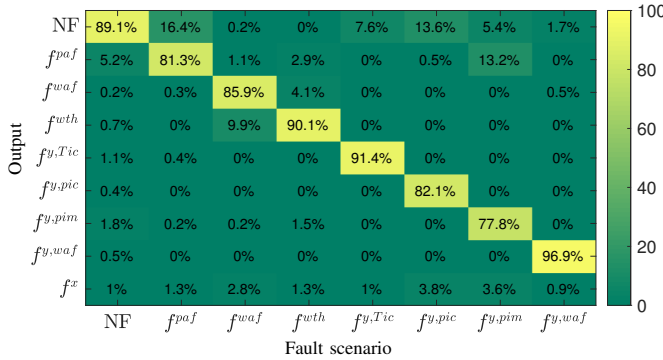


Fig. 12. Classification and ranking of a set of fault scenarios using a set of P_I -OSVM classifiers and Bayesian filtering, see Fig. 11. The rows represent the ranking of the different fault hypotheses for each fault scenario in the different columns.

The results of the unknown fault scenarios are shown in Fig. 14. Note that NF is not evaluated as an unknown fault scenario, and therefore the first column is marked with X, but it is ranked in the other fault scenarios. The unknown fault class in each fault scenario is marked with '-' since there is no P_I -OSVM model to rank that fault. In all sensor fault scenarios, i.e. $f^{y.Tic}$, $f^{y.pic}$, $f^{y.pim}$, and $f^{y.waf}$, the unknown fault class has the highest rank. The two faults f^{waf} and f^{wth} are classified as each other and f^{paf} is classified as NF, which are expected since the CAP models are overlapping, see Fig. 8.

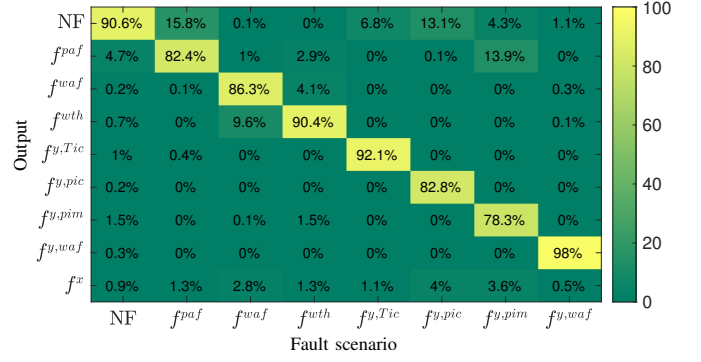


Fig. 13. Classification and ranking of a set of fault scenarios using a set of P_I -OSVM classifiers and Bayesian smoothing, see Fig. 11. There is a slight improvement compared to only using Bayesian filtering in Fig. 12.

The situation when NF gets a high rank, when a fault is present in the system, is likely when it is difficult to distinguish faults from model uncertainties and sensor noise.

One solution is to perform fault diagnosis in two steps, starting with a fault detection step followed by a fault classification step when a fault is detected. In situations where false alarms should be avoided, change detection algorithms such as cumulative sum (CUSUM) [38] can be used to reduce the false alarm rate and improve detection performance of small faults by allowing a longer time before detection. If a fault is detected with a low risk of false alarms, the following fault classification step can be performed by only considering faults without including the nominal class NF. For example, if a fault is detected in the unknown fault scenario with fault f^{paf} , see Fig. 14, and the NF fault class is removed during the Bayesian smoothing step, the ranking of f^{pim} increases from 19.6% to 82%, the unknown fault class f^x increases from 1.3% to 13%, and all the remaining fault classes remain below 2.4%. The higher ranking of f^{pim} is explained by the overlapping CAP models of the two fault classes, see Fig. 8.

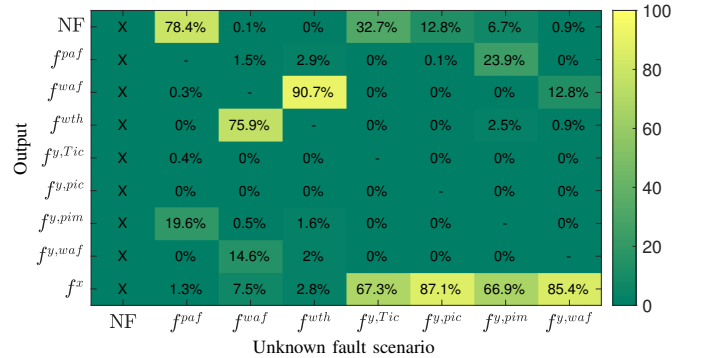


Fig. 14. Evaluation of classification of unknown fault scenarios. Training data from the selected unknown fault class are not included when training the set of P_I -OSVM models. All known fault classes are ranked using Bayesian smoothing and the evaluated fault class is marked with '-' in each scenario. Ideally, the unknown fault class f^x should have the highest rank in each column. However, some unknown faults are identified as another known fault class when the CAP models are overlapping.

The results show that the fault classification algorithm is able to handle unknown faults, but if residual data from a

new type of fault is similar to a known fault, that previously known fault class will have a higher rank. When the root cause of a detected unknown fault has been correctly identified, for example by a technician at the workshop, the fault models can be updated accordingly with the new training data, using for example incremental learning of the existing fault model or creating a new model for a new identified fault class.

VI. CONCLUDING REMARKS

Data-driven fault classification is complicated by unknown fault modes and limited training data. If multiple fault classes can explain residual data it is relevant to identify and rank the different faults instead of only selecting the most likely one, for example when supporting a technician at a workshop. The solution proposed here is to apply the principles of open set recognition which considers the problem of data classification when there are unknown fault classes and limited training data. Modeling each fault class using a P_I -OSVM classifier is used to measure the probability of inclusion that can be combined with Bayesian filtering or smoothing to improve classification performance of time-series data. An advantage of the proposed method is that it is straight forward to update and include new fault classes over time as new data are collected and labelled. Experiments using real engine data from different fault scenarios show that the proposed fault classification algorithm can identify unknown faults and that including temporal information significantly improves classification performance with respect to sample-to-sample classification.

REFERENCES

- [1] M. Blanke, M. Kinnaert, J. Lunze, M. Staroswiecki, J. Schröder, *Diagnosis and fault-tolerant control*, Vol. 2, Springer, 2006.
- [2] I. Hwang, S. Kim, Y. Kim, C. Seah, A survey of fault detection, isolation, and reconfiguration methods, *IEEE Transactions on Control Systems Technology* 18 (3) (2009) 636–653.
- [3] D. H. Stamatis, *Failure mode and effect analysis: FMEA from theory to execution*, ASQ Quality press, 2003.
- [4] D. Jung, K. Ng, E. Frisk, M. Krysander, Combining model-based diagnosis and data-driven anomaly classifiers for fault isolation, *Control Engineering Practice* 80 (2018) 146–156.
- [5] C. Sankavaram, A. Kodali, K. Pattipati, S. Singh, Incremental classifiers for data-driven fault diagnosis applied to automotive systems., *IEEE Access* 3 (2015) 407–419.
- [6] M. Atoui, A. Cohen, S. Verron, A. Kobi, A single bayesian network classifier for monitoring with unknown classes, *Engineering Applications of Artificial Intelligence* 85 (2019) 681–690.
- [7] A. Pernestål, M. Nyberg, H. Warnquist, Modeling and inference for troubleshooting with interventions applied to a heavy truck auxiliary braking system, *Engineering Applications of Artificial Intelligence* 25 (4) (2012) 705–719.
- [8] W. Scheirer, A. de Rezende Rocha, A. Sapkota, T. Boulton, Toward open set recognition, *IEEE transactions on pattern analysis and machine intelligence* 35 (7) (2013) 1757–1772.
- [9] E. Frisk, M. Krysander, Residual selection for consistency based diagnosis using machine learning models, in: *IFAC SafeProcess*, Warsaw, Poland, 2018.
- [10] A. Haghani, T. Jeansch, M. Roepke, S. X. Ding, N. Weinhold, Data-driven monitoring and validation of experiments on automotive engine test beds, *Control Engineering Practice* 54 (2016) 27–33.
- [11] A. Theissler, Detecting known and unknown faults in automotive systems using ensemble-based anomaly detection, *Knowledge-Based Systems* 123 (2017) 163–173.
- [12] H. Yan, J. Zhou, C. Pang, New types of faults detection and diagnosis using a mixed soft & hard clustering framework, in: *2016 IEEE 21st International Conference on Emerging Technologies and Factory Automation (ETFA)*, IEEE, 2016, pp. 1–6.
- [13] C. Sankavaram, B. Pattipati, A. Kodali, K. Pattipati, M. Azam, S. Kumar, M. Pecht, Model-based and data-driven prognosis of automotive and electronic systems, in: *IEEE International Conference on Automation Science and Engineering*, 2009, pp. 96–101.
- [14] C. Svärd, M. Nyberg, E. Frisk, M. Krysander, Automotive engine FDI by application of an automated model-based and data-driven design methodology, *Control Engineering Practice* 21 (4) (2013) 455–472.
- [15] J. Luo, M. Namburu, K. Pattipati, L. Qiao, S. Chigusa, Integrated model-based and data-driven diagnosis of automotive antilock braking systems, *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* 40 (2) (2010) 321–336.
- [16] D. Jung, C. Sundström, A combined data-driven and model-based residual selection algorithm for fault detection and isolation, *Transactions on Control Systems Technology* (99) (2017) 1–15.
- [17] K. Tidirri, N. Chatti, S. Verron, T. Tiplica, Bridging data-driven and model-based approaches for process fault diagnosis and health monitoring: A review of researches and future challenges, *Annual Reviews in Control* 42 (2016) 63–81.
- [18] H. Khorasgani, G. Biswas, A methodology for monitoring smart buildings with incomplete models, *Applied Soft Computing* 71 (2018) 396–406.
- [19] W. Zhang, G. Biswas, Q. Zhao, H. Zhao, W. Feng, Knowledge distilling based model compression and feature learning in fault diagnosis, *Applied Soft Computing* (2019) 105958.
- [20] K. Tidirri, T. Tiplica, N. Chatti, S. Verron, A generic framework for decision fusion in fault detection and diagnosis, *Engineering Applications of Artificial Intelligence* 71 (2018) 73–86.
- [21] I. Matei, M. Zhenirovskyy, J. de Kleer, A. Feldman, Classification-based diagnosis using synthetic data from uncertain models, in: *PHM Society Conference*, Vol. 10, 2018.
- [22] Y. Yan, P. Luh, K. Pattipati, Fault diagnosis of components and sensors in hvac air handling systems with new types of faults, *IEEE Access* 6 (2018) 21682–21696.
- [23] W. Scheirer, L. Jain, T. Boulton, Probability models for open set recognition, *IEEE transactions on pattern analysis and machine intelligence* 36 (11) (2014) 2317–2324.
- [24] E. Rudd, L. Jain, W. Scheirer, T. Boulton, The extreme value machine, *IEEE transactions on pattern analysis and machine intelligence* 40 (3) (2018) 762–768.
- [25] D. Jung, Engine fault diagnosis combining model-based residuals and data-driven classifiers, in: *IFAC International Symposium on Advances in Automotive Control*, 2019.
- [26] L. Jain, W. Scheirer, T. Boulton, Multi-class open set recognition using probability of inclusion, in: *European Conference on Computer Vision*, Springer, 2014, pp. 393–409.
- [27] G. Kitagawa, Non-gaussian state—space modeling of nonstationary time series, *Journal of the American statistical association* 82 (400) (1987) 1032–1041.
- [28] R. Domingues, M. Filippone, P. Michiardi, J. Zouaoui, A comparative evaluation of outlier detection algorithms: Experiments and analyses, *Pattern Recognition* 74 (2018) 406–421.
- [29] B. Schölkopf, R. Williamson, A. Smola, J. Shawe-Taylor, J. Platt, et al., Support vector method for novelty detection., in: *NIPS*, Vol. 12, Citeseer, 1999, pp. 582–588.
- [30] D. Tax, R. Duin, Support vector data description, *Machine learning* 54 (1) (2004) 45–66.
- [31] T. Hastie, R. Tibshirani, J. Friedman, J. Franklin, The elements of statistical learning: data mining, inference and prediction, *The Mathematical Intelligencer* 27 (2) (2005) 83–85.
- [32] D. Tax, Dtdtools, the data description toolbox for matlab, version 2.1.2 (June 2015).
- [33] J. Platt, Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods, *Advances in large margin classifiers* 10 (3) (1999) 61–74.
- [34] L. Eriksson, S. Frei, C. Onder, L. Guzzella, Control and optimization of turbo charged spark ignited engines, in: *IFAC World Congress*, 2002.
- [35] L. Eriksson, Modeling and control of turbocharged si and di engines, *OGST-Revue de l'IFP* 62 (4) (2007) 523–538.
- [36] E. Frisk, M. Krysander, D. Jung, A toolbox for analysis and design of model based diagnosis systems for large scale models, in: *IFAC World Congress*, Toulouse, France, 2017.
- [37] *Matlab 2018b statistics and machine learning toolbox*, the MathWorks, Natick, MA, USA (2018).
- [38] E. Page, Continuous inspection schemes, *Biometrika* 41 (1/2) (1954) 100–115.