

Fault Detection and Classification in Oil Wells and Production/Service Lines Using Random Forest

Matheus A. Marins^{a,*}, Bettina D. Barros^{a,b}, Ismael H. Santos^c, Daniel C. Barrionuevo^c, Ricardo E. V. Vargas^d, Thiago de M. Prego^e, Amaro A. de Lima^e, Marcello L. R. de Campos^a, Eduardo A. B. da Silva^a, Sergio L. Netto^a

^a*Electrical Engineering Program, Federal University of Rio de Janeiro, POBox 68504, Rio de Janeiro, RJ, 21947-970, Brazil*

^b*Department of Electronic Systems, Norwegian University of Science and Technology, Trondheim, Norway*

^c*Petróleo Brasileiro S.A., Cidade Universitária, Rio de Janeiro, RJ, 20211-160, Brazil*

^d*Petróleo Brasileiro S.A., Barro Vermelho, Vitória, ES, 29057-570, Brazil*

^e*CEFET-NI, Nova Iguaçu, RJ, 26041-271, Brazil*

Abstract

This paper deals with the automatic detection and classification of faulty events during the practical operation of oil and gas wells and lines. The events considered here are part of the publicly available 3W database developed by Petrobras, the Brazilian oil holding. Seven fault classes are considered, with distinct dynamics and patterns, as well as several instances of normal operation. A random forest classifier is employed with different statistical measures to identify each fault type. Three experiments are devised in order to evaluate the system performance in distinct classification scenarios. An accuracy rate of 94% indicates a successful performance for the proposed system in detecting real events. Also, the system's time of detection was on average 12% of the transient period that precedes the fault steady-state.

Keywords: fault detection and classification; oil well monitoring; abnormal event management; random forest classifier; machine learning.

1. Introduction

Recent advances in data acquisition, transmission, and storage have led many companies to develop large databases associated with their entire production and management chains. Such databases can integrate several kinds of data acquired from all sorts of sources, such as equipment sensors, economic series, human resources, and so on [1, 2]. This trend becomes a two-edged knife when the amount of data surpasses its processing capability, making it difficult to extract useful information from it. Recent developments in the

*Corresponding author.

7 fields of machine learning, computational intelligence, and data mining, however, have devised efficient algo-
8 rithms for processing large and diverse databases, capable of providing meaningful insights from underlying
9 patterns [3, 4, 5].

10 Condition-based monitoring (CBM) is a strategy that verifies the true condition of a system or equipment
11 during its continuous operation. Thus, in contrast to planned and preventive maintenance, CBM can
12 efficiently anticipate production chain problems. CBM has the following advantages [6]: it does not interrupt
13 production to evaluate equipment behavior; it promotes a safer environment for production chain workers; it
14 minimizes costs related to activity planning. According to [6], a CBM system may be structured into three
15 main phases: data acquisition, data processing, and maintenance decision making. The first phase is usually
16 the CBM bottleneck, since it requires an existing infrastructure for the acquisition of reliable data. In the
17 second phase, the data is received and processed in order to better suit the next stage [7]. The third and
18 final CBM stage is the phase in which most researchers spend their time pursuing innovative solutions [8, 9].

19 Machine learning is the most explored approach for decision-making. For instance, in this context, Yam
20 et al. [10] used recurrent neural networks to develop an intelligent predictive decision support system for
21 CBM, where Xavier and Seixas [11] applied a similar algorithm to analyze a chemical process. Widodo and
22 Yang [12] and Helmy et al. [13] also used machine learning approaches based on support vector machines to
23 model the decision-making process.

24 This work proposes a CBM system for oil and gas (O&G) wells and production/service lines that acts as a
25 support tool to decision-making systems. The proposed system attempts to detect and identify an anomalous
26 behavior as early as possible, so the operator, given the fault inherent criticality, can intervene accordingly
27 to reduce the associated production losses. The work has been developed using a database developed by
28 Petrobras comprising about 2000 events. They describe the normal operation (Class 0 in this work) and
29 eight different types of fault (Classes 1 to 8) in well operation [14]. Our CBM system is an ensemble of
30 machine learning techniques, from the pre-processing phase until the decision-making. The proposed system
31 starts by extracting nine different statistical features from each available tag, and then performs a principal
32 component analysis [3] to reduce the problem dimension as well as the associated computational complexity.
33 This is followed by a random forest classifier, that is used to detect and classify the faulty conditions, aided
34 by a Bayesian approach [15] employed to tune the classifier hyperparameters.

35 Three different classification scenarios are devised: fault detector, single-fault detector/classifier, and
36 multi-fault classifier. Results in each experiment reach up to 90% of system accuracy, indicating its ability
37 to detect and classify most of the fault types properly. In order to introduce its technical contributions, this
38 paper is organized as follows: Section 2 describes the general problem of detecting faulty behaviors in O&G
39 wells and production/service lines; Section 3 introduces the proposed CBM system and Section 4 details
40 the database employed in this work, the so-called 3W database; Sections 5 and 6 present the experimental
41 methodology and results, respectively, detailing and analyzing the performance of the proposed system in

42 the problem at hand. Finally, Section 7 closes the paper emphasizing its main contributions.

43 2. Problem Description

44 The production of O&G from underground reservoirs involves chemical and mechanical processes that
45 affect well drilling and operation. Many of these processes may eventually cause a problem with the well,
46 resulting in a decrease in production or in equipment failure. The majority of serious problems can be
47 avoided or postponed by preventive maintenance techniques, or recognized at an early stage by means of
48 regular analysis of production rates, fluids, and the mechanical condition of the well. Such practices can
49 prevent expensive workovers that may be necessary to restore well production and can also prevent total
50 well losses.

51 2.1. O&G General Fault Assessment

52 Since 2010, Petrobras has been developing a large database aimed at describing all of its O&G losses
53 within its operational unit in Rio de Janeiro (OU-Rio). The so-called loss integrated management (LIM)
54 platform congregates several complementary databases describing a production loss by a series of up to 85
55 different features such as: initial/final loss date, platform, affected equipment, equipment operator, related
56 sensor tags, original cause, secondary cause, estimated loss, required initial/secondary actions, subsequent
57 activities, and so on [16, 17]. Figure 1, for instance, shows the relative cumulative volume loss and number of
58 failures between 2014 and 2017 of Petrobras OU-Rio which comprises 298 production/injection wells. From
59 this figure, one can observe the significant amount of loss (23.8%) caused by reservoir- and well-related issues.
60 A breakdown of these losses is depicted in Figure 2, indicating the main origins of the faults associated with
61 such losses.

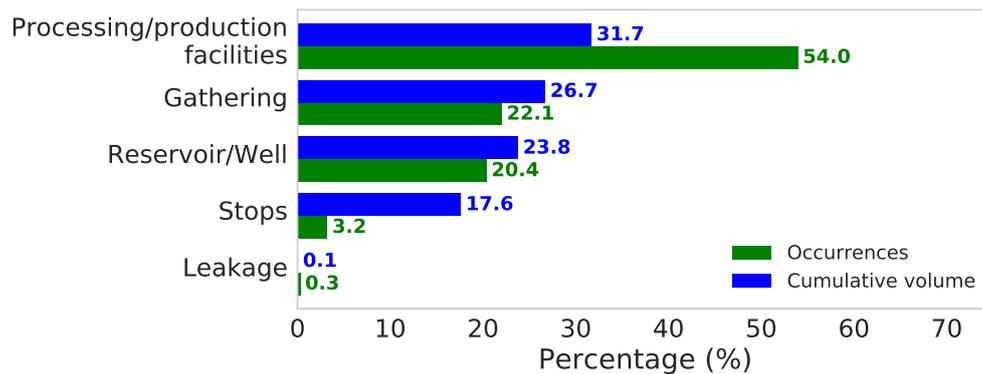


Figure 1: Breakdown of cumulative oil-volume loss (blue bars) and corresponding number of failures (green bars) between 2014 and 2017 in Petrobras OU-Rio for different production systems.

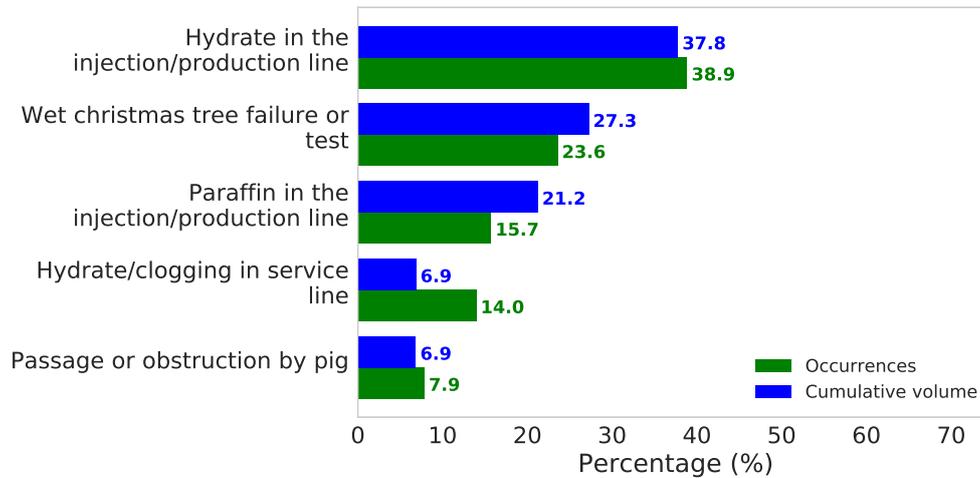


Figure 2: Breakdown of cumulative oil-volume loss (blue bars) and corresponding number of failures (green bars) between 2014 and 2017 in Petrobras OU-Rio for different fault causes.

2.2. Well/Line-Related Fault Description

In mid-2017, Petrobras conceived a project, entitled monitoring of specialized alarms, to develop a new automated system for detecting and classifying eight types of undesirable events in offshore naturally flowing wells (see the 3W database described in Section 4). The selected types of events are [14]:

Class 1 – Abrupt Increase of Basic Sediment and Water (BSW): The BSW is the ratio between water and sediment flow rate and the liquid flow rate, both measured under normal temperature and pressure, providing an insight on the amount of water in the produced oil. During the life cycle of a well, it is expected that the BSW increases due to water production from either the natural reservoir aquifer or artificial injection. An abrupt increase of BSW, however, is often associated with multiple problems related to flow assurance, lower oil production, oil lifting, incrustation, and so on. Early identification of this type of event helps to avoid these undesired conditions.

Class 2 – Spurious Closure of the Downhole Safety Valve (DHSV): The DHSV is installed in the production tubing of wells to ensure their closure in emergency scenarios or physical disconnection. Sometimes the closure function fails in a spurious manner without any indication on the surface. Production losses can be avoided if this spurious closure is detected as precociously as possible, allowing corrective operational procedures at an early stage.

Class 3 – Severe Slugging: Stratified gas-liquid flow, which occurs in scenarios of low liquid and gas flow rate, and a declined production line tend to cause liquid accumulation at the bottom of the riser, which blocks the gas flow until sufficient upstream pressure causes a surge of liquid and gas. After this sudden surge, some of the liquid in the riser returns to the base blocking the flow once again thus restarting the cycle. This transient cyclic phenomenon is called severe slugging and, depending on its periodicity and intensity, may become critical as it leads to stress or even damage to well equipment. Early detection of

84 this type of event allows preemptive actions to reverse the situation before it becomes critical.

85 **Class 4 – Flow Instability:** This type of undesirable event is also characterized by spurious surges of
86 liquid and gas, as is the case of severe slugging, with the difference that it is less intense and does not involve
87 the complete cycle of liquid blockage followed by a gas surge. If not dealt with accordingly, flow instability
88 can evolve to severe slugging.

89 **Class 5 – Rapid Productivity Loss:** Productivity of a naturally flowing well (as opposed to fluid
90 injected wells) depends on several conditions. When the system energy is less than the minimum necessary
91 to overcome energy loss, for instance, the oil flow slows or even stops. Therefore, early detection of this
92 undesired condition reduces production losses.

93 **Class 6 – Quick Restriction in the Production Choke (PCK):** The PCK installed in the produc-
94 tion unit is responsible for controlling the well from the surface. Manual operation of this type of valve can
95 lead to unwanted quick restrictions, affecting the oil production directly.

96 **Class 7 – Scaling in PCK:** Another PCK-related undesirable event is the scaling due to inorganic
97 deposits along time, which can severely reduce oil production. Early identification of this type of event
98 is also desired, as special actions can be taken, such as injection of scale inhibitors, to avoid additional
99 production losses.

100 **Class 8 – Hydrate in Production Line:** As shown in Figure 2, the presence of hydrate in wells
101 and in production/injection lines is one of the biggest problems in the O&G industry, including Petrobras.
102 Hydrates are crystalline compounds, resembling ice in appearance, formed by the reaction of natural gas
103 to water. Early detection of this undesirable event means avoiding production losses for long periods, as
104 unblocking production lines is costly and sometimes a long process.

105 During the well operation, it is possible that these faults interact with each other. For example, two
106 or more faults can coincide, for instance, a Spurious Closure of DHSV (Class 2) suddenly occurs during a
107 Scaling in the PCK (Class 7), which is a very slow event. In this case, the scaling should be detected before
108 the closure occurs; otherwise, it would be impossible to detect it until the production restarts.

109 It is also possible that a fault triggers another fault from a different class. For instance, Flow Insta-
110 bility may precede Severe Slugging because these two faults often occur in scenarios of low liquid and gas
111 flowrates and wavy flowlines. When the well starts to oscillate, depending on the line geometry and the flow
112 characteristics, this instability can evolve to a more dangerous scenario of severe slugging.

113 Also, Severe Slugging (Class 3) and Flow Instability (Class 4) can be caused by faults related to partial
114 flow area blockage, such as Rapid Production Loss (Class 5), Scaling and Quick Restriction in PCK (Class
115 6), and Hydrate in Production Line (Class 8). All these events may cause a reduction of the production
116 flowrate without completely stopping the well production, leading to slugging and instability.

117 In this paper, a CBM system is proposed for automatically detecting and classifying the above event
118 types using machine learning techniques. The data-driven approach learns underlying patterns during faulty

119 well operations and identifies such anomalous behaviors in subsequent system operations, as detailed in the
 120 following section.

121 3. System Overview

122 The proposed system is represented in Figure 3, which shows the data flow starting from its raw version
 123 all the way down to the system classification output, passing through the feature extraction, data trans-
 124 formation, and classification modeling stages. The first two blocks extract information in a compact form,
 125 whereas the last block associates the classification labels with the related characteristics according to the
 126 task at hand. Details of each system stage are provided in the three ensuing subsections.

127 Every system block was implemented using Python, which has become one of the standard programming
 128 languages for machine learning algorithms, due to its readability, growing community, and available libraries.
 129 In particular, in this work, we use two of these major libraries: Pandas [18], for data analysis, and Scikit-
 130 Learn [19], for handling the basic algorithms.



Figure 3: Block diagram of the proposed CBM system for O&G wells and production/injection lines.

131 3.1. Feature Extraction

132 Feature extraction is commonly concerned with highlighting important information to help the classifi-
 133 cation task. In the proposed system we process the raw input data from the available tags, all collected at
 134 a rate of one sample per second, and extract n_f features from each N -sample data window for each tag.
 135 Given an N -sample tag window $\mathbf{x} = [x_1, x_2, \dots, x_N]^T$, the $n_f = 9$ extracted features in the proposed system
 136 are the following statistical measurements:

- 137 • **Mean value:**

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i; \quad (1)$$

- 138 • **Standard deviation:**

$$\sigma = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)^2}, \quad (2)$$

139 which provides information on how the data is spread around its mean;

- 140 • **Skewness:** The third standardized moment of a random variable [20],

$$s_k = \frac{E[(X - \mu)^3]}{(E[(X - \mu)^2])^{\frac{3}{2}}}, \quad (3)$$

141 quantifies the asymmetry of the given data distribution: $s_k = 0$ stands for a symmetric distribution,
 142 where the mean is equal to the median;

143 • **Kurtosis:** The fourth standardized moment of a random variable [21],

$$k = \frac{E[(X - \mu)^4]}{(E[(X - \mu)^2])^2} - 3, \quad (4)$$

144 measures the tails of the distribution, which are related to the presence of outliers. Positive and
 145 negative kurtosis indicate heavy-tailed and light-tailed distributions, respectively.

146 • **5-number summary:** The last five features are given by the minimum, maximum, median, and first
 147 and third quartiles values of the data window [22]. The minimum and maximum values provide the
 148 range of the data distribution within a given window, the median is the center of the distribution, and
 149 the first and third quartiles show how the data is spread within its extremes.

150 In summary, the feature-extraction block receives an $\mathbb{R}^{M \times n}$ data matrix, with M time samples of n input
 151 tags, and transforms it into an $\mathbb{R}^{m \times n_f}$ matrix, where $m = \lfloor \frac{M-N}{s} \rfloor$ is the total number of N -sample data
 152 windows, s denotes the shift in samples between two consecutive windows, and $n_f = 9$ is the number of
 153 features extracted from each data window.

154 3.2. Data Transformation

155 The concept of data transformation is to represent data features in a more appropriate space. In our
 156 context, its purpose is twofold. First, it avoids range divergence among the extracted features. Second, it
 157 mitigates the “curse of dimensionality” problem, which means that although high dimensionality is expected
 158 to provide more information from the data, it also affects negatively its discriminative capability [23]. The
 159 first goal is achieved by performing z -score normalization on the data, and the second one by employing
 160 principal component analysis (PCA) on the result of the z -score normalization.

161 The z -score or standardized normalization [24] tends to equalize the features ranges and is performed as
 162 follows:

$$Z^i = \frac{X_f^i - E[X_f^i]}{\sigma_f}, \quad (5)$$

163 where X_f^i is a random variable representing one of our data features, σ_f is its standard deviation, and Z^i is
 164 the normalized version of X_f^i .

165 PCA [3] is a technique for reducing the data dimensionality while keeping most of its energy, that is
 166 in general equivalent to keep most of its representative information. Assuming a centered data matrix
 167 $\mathbf{X}_f = [\mathbf{x}_f^1, \dots, \mathbf{x}_f^m]^T \in \mathbb{R}^{m \times n_f}$, with a sample covariance matrix $\mathbf{S} = (m - 1)^{-1} \mathbf{X}_f \mathbf{X}_f^T$, the PCA formulation
 168 is given by

$$\mathbf{V}\mathbf{S} = \mathbf{\Lambda}\mathbf{V}, \quad (6)$$

169 where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m)$ is the diagonal matrix of eigenvalues λ_i of \mathbf{S} and $\mathbf{V} = [\mathbf{v}^1, \dots, \mathbf{v}^m]^T$ is the
170 orthonormal matrix of eigenvectors \mathbf{v}^i , also known as principal components, of \mathbf{S} . The PCA transformation
171 $\mathbf{X}_f^m = \mathbf{V}\mathbf{X}_f$ performs a data projection into a new set of coordinates that are sorted in decreasing order
172 of variance directions (i.e., \mathbf{v}^1 and \mathbf{v}^m represent the maximum-variance and minimum-variance directions,
173 respectively).

174 The dimensionality reduction consists in selecting only a small number $p < m$ of principal components
175 $\mathbf{V}^p = [\mathbf{v}^1, \dots, \mathbf{v}^p]^T$, which retain a large fraction (99%, for instance) of the original data information/energy,
176 thus reducing the feature matrix dimensions from $\mathbb{R}^{m \times n_f}$ to $\mathbb{R}^{p \times n_f}$.

177 3.3. Classification Modeling

178 In the past few years, deep neural networks (DNN) have been applied in several O&G-related problems,
179 including failure classification during drilling operations [25]. However, when compared to other ML algo-
180 rithms, DNNs have several disadvantages such as the requirement of a large and adequately labeled dataset,
181 the vast choices for the hyperparameter tuning (number of layers, number of neurons in each layer, type of
182 nonlinear activation function, and so on), the complexity of its training procedure (due to a large number
183 of internal coefficients), and its sensitivity to outliers or missing input data.

184 In this paper, we adopted the random forest [26] algorithm, which is a supervised machine-learning
185 approach suitable for both binary and multiclass problems such as the ones considered here. The random
186 forest technique has already been applied to many distinct problems, such as gene selection [27], remote
187 sensing [28], prediction of proteins [29], and so on. The main characteristics of this algorithm are:

- 188 • it is robust to noise and outliers;
- 189 • it is faster than bagging and boosting methods [3, 5];
- 190 • it can provide useful error information (strength, correlation, and importance of the variable);
- 191 • it is simple to deploy (as it has a small number of hyperparameters to be tuned);
- 192 • it is easy to parallelize;
- 193 • it can handle missing data.

194 The random forest algorithm is an ensemble of decision trees, known as weak learners, for having low
195 computational cost and low discriminative capabilities. Training a random forest classifier is equivalent to
196 training several independent decision trees. When training each of these trees, distinct subsets of the input
197 data and the extracted features are randomly drawn, so that each tree learns from a different partition of
198 the data, as depicted in Figure 4. Properly combined, these trees can generate strong classifiers using the
199 idea of wisdom of the masses [3, 4, 5]. After the learning procedure a new input data sample is labeled as
200 the class that the majority of the decision tree classifiers voted for.

201 The vote distribution of all trees within a random forest can be interpreted as a probabilistic distribution
 202 for the system output. For example, if we use 100 trees to analyze a given input sample, which 80 trees
 203 estimate as Class 0 and the rest as Class 1, we can say that the classifier has 80% and 20% of certainty that
 204 this sample belongs to Classes 0 or 1, respectively. This property can be used to establish a threshold τ_i
 205 each output class, with a given sample only belonging to that class if the certainty is above the corresponding
 206 threshold. In the context of decision-making systems, these classification thresholds allow one to balance the
 207 number of false positives and undetected faults according to the system priority: increasing the threshold
 208 values lowers the number of fault misclassifications at the cost of a reduction in the fault detection rates.

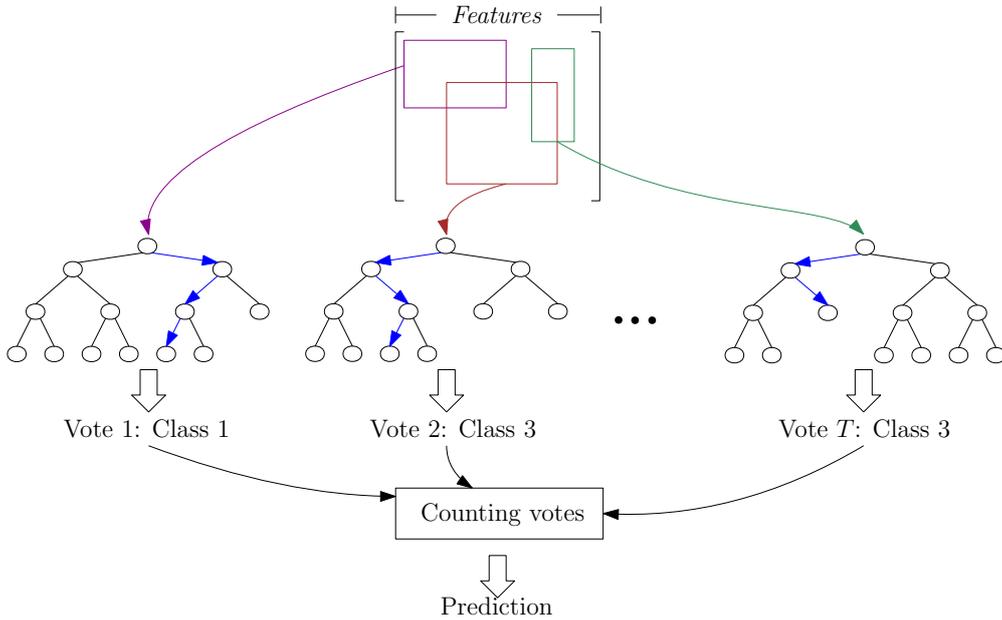


Figure 4: General schematic diagram of the random forest classifier.

209 4. The 3W Database

210 The development of the proposed machine-learning CBM system was based on the 3W dataset introduced
 211 in [14] (see also Sections 1 and 2). This database is composed by approximately 2000 operational events
 212 representing different states of the well, varying from normal operation (Class 0) to the eight distinct faults
 213 (Classes 1 to 8) described in Section 2. Each event is a time-series data composed by $n = 8$ tags acquired
 214 by 8 different sensors, chosen according to their availability and relevance to the faults at hand, as given in
 215 Table 1.

216 Three event types are included in the 3W database: real, simulated, and sketched¹. Real events are
 217 described by sensor data acquired through the PI System [30] during real well operations. Signals of

¹In [14], the authors named the sketched events as hand-drawn events.

Name	Description	Unit
P-PDG	Pressure at permanent downhole gauge (PDG)	Pa
P-TPT	Pressure at temperature/pressure transducer (TPT)	Pa
T-TPT	Temperature at temperature/pressure transducer (TPT)	°C
P-MON-CKP	Pressure upstream of production choke (CKP)	Pa
T-JUS-CKP	Temperature downstream of production choke (CKP)	°C
P-JUS-CKGL	Pressure downstream of gas lift choke (CKGL)	Pa
T-JUS-CKGL	Temperature downstream of gas lift choke (CKGL)	°C
QGL	Gas lift flow rate	m ³ /s

Table 1: List of tags in 3W database, including tag names, descriptions, and measuring units.

218 simulated origin were obtained through the OLGA system [31], vastly used by the industry for dynamic
219 multiphase flow simulation. Sketched signals were created through a tool designed by 3W database creators,
220 which uses expert knowledge to sketch the profile of a particular unwanted event. Table 2 contains the
221 quantitative description of 3W dataset per event type.

Class	Description	Real	Simulated	Sketched	Total
0	Normal	597	0	0	597
1	Abrupt BSW Increase	5	114	10	129
2	Spurious DHSV Closure	22	16	0	38
3	Severe Slugging	32	74	0	106
4	Flow Instability	344	0	0	344
5	Rapid Productivity Loss	12	439	0	451
6	Quick PCK Restriction	6	215	0	221
7	Scaling in PCK	4	0	10	14
8	Hydrate in Prod. Line	3	81	0	84
	Total	1025	939	20	1984

Table 2: Quantitative description of 3W database per event type.

222 When building the 3W database, besides the faulty periods, the authors have indicated for each event
223 instance a period when the samples are not faulty (referred to as normal) and a period of fault transient
224 before the actual steady-state fault consolidation, as depicted in Figure 5. Such annotation procedure allows
225 one to detect a given fault during its initial transient stage. This is the behavior intended in the proposed
226 system, in order to minimize maintenance costs and production losses.

227 5. Experimental Methodology

228 In this work, three experiments were devised in order to evaluate and understand the capability of the
229 proposed CBM system:

230 **Experiment 1 – One-class classifier:** This scheme consists of a single classifier to discriminate only
231 between normal and faulty operations. Therefore, in this case, all faults are combined into a unique class
232 which is compared against the normal-operation class.

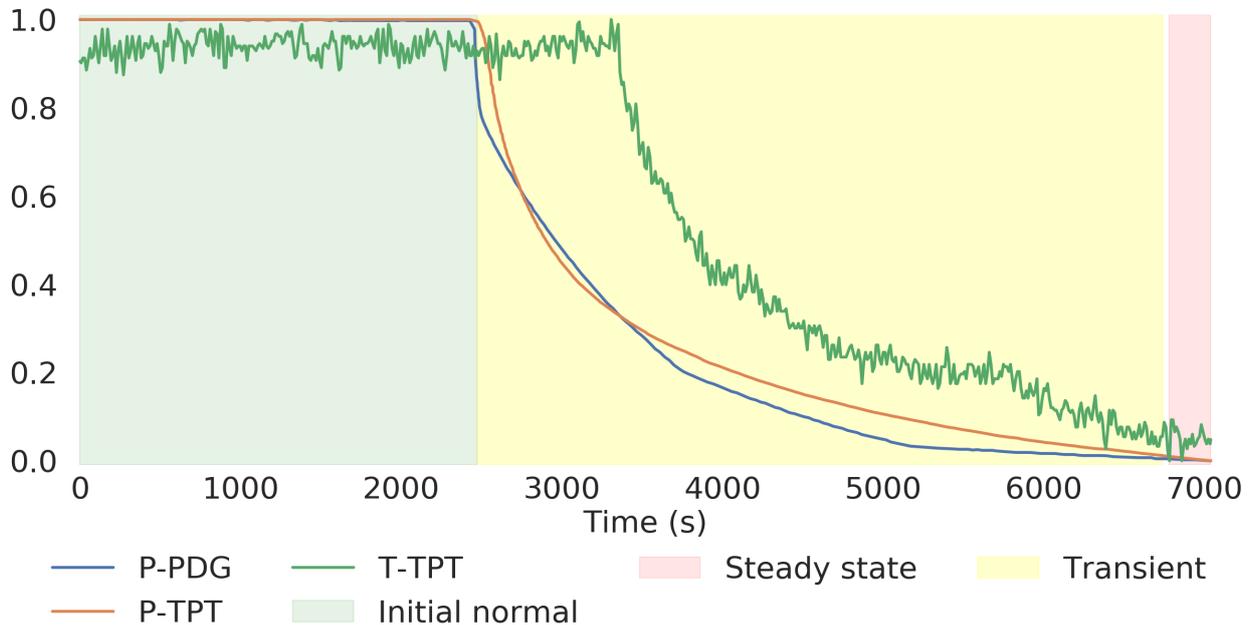


Figure 5: Example of (normalized) tag values of a Class 2 event where the background color indicates the normal (green), transient (yellow), and faulty (red) stages.

233 **Experiment 2 – Multiple binary classifiers:** This strategy consists in designing several classifiers,
 234 each one specialized in discriminating an individual fault against the normal-operation mode. In this strategy,
 235 one can infer which faults are the hardest to be identified, as opposed to the previous scheme which analyzes
 236 the faulty conditions altogether in a single class.

237 **Experiment 3 – Single multiclass classifier:** This scheme employs a single system to identify all
 238 different classes individually. In this case, each class (normal operation or any specific fault) is discriminated
 239 against all the remaining classes, thus providing more information to the system operator but posing a much
 240 harder problem to the machine-learning algorithm.

241 As described in Section 4, each annotated event in the 3W dataset has three phases: normal, transient,
 242 and steady state. Our main research goal was building a CBM system capable of anticipating a fault as
 243 much as possible, in order to give as much time as possible for the operator to intervene and minimize the
 244 losses. For this reason, in all experiments we have opted to use only the normal and transient phases of
 245 the event to train the random forest classifier. Also, for a more realistic scenario, the sketched instances
 246 were ignored and Class 7 was disregarded as it is under-represented in the 3W database. After all these
 247 considerations, the 3W database was split into the training and test sets, following a 70%/30% ratio, as
 248 indicated in Table 3.

Class	Description	Real	Simulated	Total
0	Normal	468(129)	0	468(129)
1	Abrupt BSW Increase	3(2)	78(36)	81(38)
2	Spurious DHSV Closure	15(7)	11(5)	26(12)
3	Severe Slugging	22(10)	55(19)	77(29)
4	Flow Instability	260(84)	0	260(84)
5	Rapid Productivity Loss	8(4)	340(99)	348(103)
6	Quick PCK Restriction	4(2)	170(45)	174(47)
8	Hydrate in Prod. Line	0(3)	56(25)	56(28)
Total		780(241)	710(229)	1490(470)

Table 3: Number of instances in both training and test (between parenthesis) sets employed in the development of the proposed CBM system. Notice the absence of any Class 7 and sketched events.

249 5.1. Training, Validation, and Test

250 Building a system with a wide variety of hyperparameters demands a robust and reliable training routine.
251 In this work, we performed a cross-validation [32] using $k = 5$ folds. In order to avoid contamination during
252 this procedure, all samples of a given event belong to the same fold. Also, we kept the same class proportions
253 in the training set throughout every cross-validation iteration.

254 In our experiments, we have considered a hyperparameter to evaluate the amount of balance between
255 normal (including n_0 samples from Class 0 and n_N samples from the initial normal phase in the faulty
256 instances) and faulty samples. We refer to this parameter as b , and, given n_0 and n_N , we have $b = n_0/n_N$.
257 So, the total number of normal and faulty samples used during training becomes

$$n_{\text{normal}} = n_{\text{faulty}} = (b + 1)n_N. \quad (7)$$

258 Even though b is not mandatory when training a random forest classifier, incorporating it to our hyper-
259 parameter search enables us to increase the performance on the samples from the initial normal phase in
260 the faulty instances. Increasing b makes n_N less relevant in comparison to n_{faulty} , feeding the model with
261 less information about those n_N samples.

262 Many hyperparameters have a direct impact on the system architecture, and thus can influence its overall
263 performance. In the proposed system, for instance, we considered, for the PCA stage, the minimum number
264 of components that guarantee a 99% threshold for the accumulated energy. A fixed sliding window step of
265 $s = 50$ samples was also employed.

266 The strategy applied to assess the best set of the remaining hyperparameters was a combination of the
267 traditional grid search and a Bayesian optimization method. The Bayesian method applied here [15] is a
268 non-convex optimization algorithm that samples the objective function according to a Gaussian process.
269 In each optimization iteration, the algorithm considers all past observations to evaluate which parameter
270 regions are worth exploring, thus narrowing the hyperparameter search space following a probabilistic-based
271 strategy.

272 For choosing the sliding window size and the balance ratio, we performed grid search using the following
 273 sets of values: $M = \{100, 200, 300, \dots, 900, 1000\}$ and $b = \{1, 2, 3, 4, 5\}$. As for the number of trees N_t and
 274 the maximum depth M_d of each tree in the random forest algorithm, we applied 50 iterations of a non-
 275 convex optimization algorithm based on Bayesian sampling of the objective function, as described in [15].
 276 The following parameter ranges have been used: $50 \leq N_t \leq 250$ and $5 \leq M_d \leq 70$.

277 5.2. Metrics Employed for Performance Assessment

278 In this work we evaluate the classifier models through three different metrics:

- 279 • Accuracy (ACC): considering all normal and faulty samples, the ACC can be computed as:

$$\text{ACC} = \frac{TP + TN}{n_{\text{total}}}, \quad (8)$$

280 where TP and TN are respectively the number of true positives and true negatives, in samples, and
 281 n_{total} is the overall number of samples;

- 282 • Faulty-normal accuracy (FNACC): is the accuracy computed when considering only the normal sam-
 283 ples preceding a faulty instance;

- 284 • Real faulty-normal accuracy (RFNACC): is the accuracy computed when considering only the normal
 285 samples preceding a real faulty instance.

286 While the standard ACC addresses the overall model performance, it does not evaluate its capability of
 287 discriminating between the transient phases in the faulty events and the normal samples preceding them.
 288 Therefore, in order to assess the system false-alarm rate, the FNACC and RFNACC are more indicated,
 289 with the later proving insight only of the performance for real faulty events.

290 6. Experimental Results and Discussion

291 In this section, we present the results for all three classification strategies described in Section 5. This
 292 provides a comprehensive assessment of the final system performance, under different fault-detection sce-
 293 narios.

294 6.1. Experiment 1

295 After evaluating 2500 different binary models (following 50 runs of the Bayesian optimizer, for ten
 296 different values of window size M and five values of balance ratio b), the best classifiers according to the
 297 three distinct evaluation metrics are shown in Table 4, along their average validation performance.

298 For the evaluated configurations, we may state that larger values of b are associated to larger values of
 299 ACC and smaller values of RFACCC, as depicted in Figure 6. This happens because a larger b leads to an
 300 increased Class 0 sample representation in comparison to the normal-phase samples preceding faults in the

Model	M	b	N_t	M_d	ACC	FNACC	RFNACC
1.1	700	5	179	70	0.986	0.935	0.660
1.2	600	1	86	53	0.979	0.965	0.733
1.3	500	2	118	38	0.983	0.961	0.773

Table 4: Characterization of best classifiers in Experiment 1, each one according to a specific accuracy metric, along their validation performances: window length M , normal-class balance ratio b , random forest number of trees N_t , and maximum tree depth M_d .

301 training datasets. From Figure 6, one may also notice that the window length M and the number of trees
302 N_t do not have a strong impact on the final results, whereas the maximum tree depth M_d may hinder the
303 final performance in this experiment if one chooses $M_d < 15$.

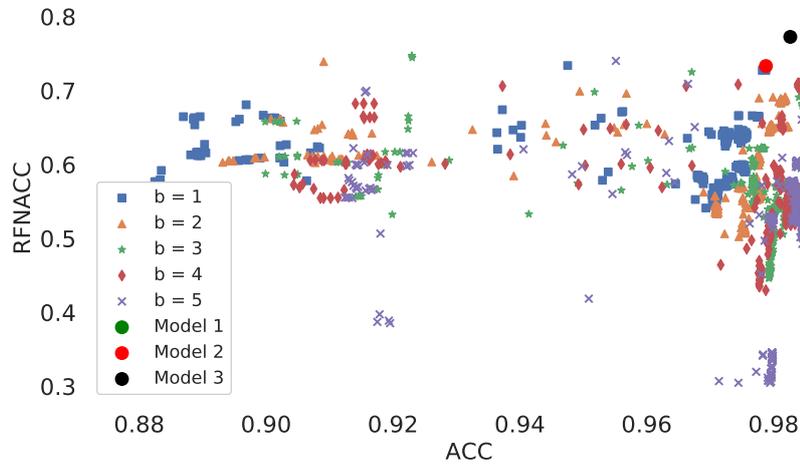


Figure 6: ACC \times RFACC relationship for all 2500 classifier configurations evaluated in Experiment 1, indicating that larger balance ratios b correspond to larger ACC and smaller RFACC values.

304 From a practical perspective, analyzing the model performance in the real instances is far more critical
305 than analyzing its performance in the simulated ones. Moreover, the model must avoid false positives, as
306 each time this happens, the operator loses confidence in using the system. For these reasons, we chose to
307 analyze Model 1.3 test performance. The breakdown of the test result is shown in Table 5, including its
308 performance in classifying the steady-state phase of the faults (which were not used in training).

Event	Class 0	Initial Normal	Transient	Steady-State	Overall
Simulated	–	0.990	0.995	0.958	0.971
Real	0.995	0.673	0.822	0.972	0.971
Overall	0.995	0.938	0.982	0.960	0.971

Table 5: Model 1.3 test accuracy results for each separate sample type, including simulated/real events. The empty entry indicates the absence of data for that particular event type.

309 Although Model 1.3 was able to reach a very large overall accuracy of ACC = 0.971, Table 5 results
310 show its difficulty in detecting properly the normal samples preceding faults in real events, what may lead

311 to an undesired large number of false alarms. One of the possible causes for this poor behavior is the system
 312 difficulty in modeling the distinct dynamics of all different faults combined into a single class. The next
 313 experiment avoids this issue by designing individual classifiers for each separate fault.

314 6.2. Experiment 2

315 In this experiment, we select the three configuration models chosen in the previous experiment (that
 316 optimize each of the three distinct performance measures), as given in Table 4, and retrain them under the
 317 new circumstances. In this case, we use a different classifier to distinguish between the normal and faulty
 318 operations for each fault type. Table 6 summarizes all the test results for each model and for each fault
 319 type (Classes 1 to 8, except Class 7). As one can notice, in general, Models 1, 2 and 3 yield larger ACC,
 320 FNACC, and RFNACC values, as each model tends to prioritize each of these metrics, respectively.

Fault Class	Model	ACC	FNACC	RFNACC
1	1.1_1	0.992	0.932	0.170
1	1.2_1	0.994	0.975	0.704
1	1.3_1	0.994	0.982	0.787
2	1.1_2	0.999	0.994	0.992
2	1.2_2	0.998	0.985	0.970
2	1.3_2	0.998	0.981	0.963
3	1.1_3	1.000	–	–
3	1.2_3	1.000	–	–
3	1.3_3	1.000	–	–
4	1.1_4	0.990	–	–
4	1.2_4	0.989	–	–
4	1.3_4	0.989	–	–
5	1.1_5	0.984	–	0.503
5	1.2_5	0.962	–	0.965
5	1.3_5	0.968	–	0.875
6	1.1_6	0.966	0.976	0.999
6	1.2_6	0.963	0.967	0.923
6	1.3_6	0.962	0.977	1.000
8	1.1_8	0.969	0.994	1.000
8	1.2_8	0.968	0.998	1.000
8	1.3_8	0.969	0.992	1.000

Table 6: Experiment 2 test results for each classifier model and for each fault type (Classes 1 to 8, except Class 7). Empty entries indicate the absence of data for that particular event type. The models are named as 1.X_Y, where X corresponds to the model chosen in Experiment 1 and Y to the Fault Class index.

321 These results show that the proposed system can properly detect and classify all faults, particularly the
 322 ones of Classes 2, 3, 4, 6, and 8, where the accuracy levels reached 0.9 or higher for every metric. Even
 323 Classes 1 and 5 yielded good accuracy values but presented lower RFNACC when compared to other classes.

324 This experiment shows that the developed system can identify properly all individual faults against the
 325 normal operation, even when using the same model parameters for all classes. As pointed out in Section 2,

326 different faults can occur at the same time or even evolve to other fault types. Therefore, in the next
 327 experiment, we tackle these scenarios altogether using a multiclass strategy.

328 6.3. Experiment 3

329 In this scenario, a new hyperparameter search for the multiclass system was performed considering the
 330 same ranges explored in Experiment 1. Table 7 brings the results for the three system models considering
 331 the different accuracy metrics as before.

Model	M	b	N_t	M_d	ACC	FNACC	RFNACC
3.1	500	5	102	24	0.973	0.936	0.515
3.2	400	1	86	27	0.962	0.963	0.715
3.3	800	2	238	30	0.970	0.962	0.719

Table 7: Characterization of best classifiers in Experiment 3, each one according to a specific accuracy metric, along their validation performances: window length M , normal-class balance ration b , random forest number of trees N_t , and maximum tree depth M_d .

332 When comparing Tables 4 and 7, one readily notices the similar performances achieved in each case,
 333 despite the different classification problem. This can be explained by the difficulty posed in Experiment 1
 334 for grouping all fault types into a single class, not respecting their distinct natures and dynamics, what can
 335 be inferred from the results obtained in Experiment 2.

336 A breakdown of all classification test results obtained by Model 3.3 is provided in Table 8. These results
 337 correspond to an overall system accuracy of 0.94 and excellent classification rates specially for Classes 3, 4,
 338 5, and 6.

339 We can also notice that classes with less samples have worst performances, as less data hinders the
 340 learning process during the system training. This culminates in Class 8, which had no real samples in the
 341 training set and presented the lowest test accuracy, indicating that the simulated instances do not emulate
 342 properly the real cases. To verify this hypothesis, an experiment was conducted where Model 3.3 was
 343 retrained with the number of Class-4 events reduced to only 10% of the available instances. In that scenario,
 344 the overall test accuracy dropped to 92.3% and the accuracy within the Class 4 dropped from 95.5% to only
 345 51.4%, as explained above.

346 Also, an overall confusion matrix is shown in Figure 7, which shows that the model successfully identifies
 347 almost all fault cases. The misclassification of type-8 into type-3 faults, as seen in this matrix, can be
 348 explained by the fact that severe slugging (Class 3) often results from hydrate in the production line (Class
 349 8), which is an example of flow area blockage, as mentioned in the end of Section 2

350 These error patterns may be mitigated in a real scenario by providing additional side information to the
 351 system operators, thus enabling them to sort out the system output based on their own practical knowledge.

Event	Type	Initial Normal	Transient	Steady-State	Overall
Class 0	Real	–	–	–	0.994
Class 1	Real	0.719	0.454	0.000	0.508
	Simulated	1.000	0.996	0.944	0.978
Class 2	Real	0.982	0.882	0.817	0.888
	Simulated	0.978	0.849	0.035	0.234
Class 3	Real	–	–	0.791	0.791
	Simulated	–	–	0.956	0.956
Class 4	Real	–	–	0.954	0.954
	Simulated	–	–	–	–
Class 5	Real	0.443	0.953	–	0.833
	Simulated	–	0.997	0.940	0.949
Class 6	Real	0.832	0.153	0.146	0.710
	Simulated	0.994	0.983	0.873	0.908
Class 8	Real	0.000	0.000	0.000	0.000
	Simulated	0.998	0.985	1.000	0.989
Overall	Real	0.615	0.511	0.901	0.930
	Simulated	0.998	0.992	0.918	0.944
Overall	–	0.934	0.955	0.916	0.940

Table 8: Model 3.3 test accuracy results for each separate sample type, including simulated/real events. Empty entries indicate the absence of data for that particular event type.

352 6.4. Event-Based Assessment

353 Regarding the system operating in a real scenario, its main goal is detecting an event as soon as possible
354 despite misclassifying some time samples. Therefore, the goal of this subsection is to analyze the performance
355 of Model 3.3, which focuses on real events, at the event level.

356 We consider that an event is correctly detected and classified if the model obtains an accuracy rate higher
357 than a given threshold on the samples of that particular event. In this paper, we have opted for a conservative
358 accuracy rate of 0.9 at the sample level to consider an event properly classified. In Table 9, we present the
359 classification results for each fault type, discriminating the number of correctly classified events, for each
360 event phase. Even though the classifier was not trained with steady-state samples, we used Model 3.3 to
361 predict it, considering the label of the same fault transient label. The column “All” represents the scenario
362 where the model detected a fault during its transient or during the subsequent steady-state period, and
363 correctly detect the initial normal phase. These event-level results show a quite satisfactory performance
364 for the proposed system, with low false-alarm rates and high true-positive identification/classification rates.

365 In Figure 8, we have an example of a sample-level classification (red crosses), where the ‘0’ and ‘1’
366 outputs correspond to normal and faulty states, respectively. Two new time-intervals (in seconds) are also
367 defined in this figure: the time of detection, t_1 , between the beginning of the transient phase and the fault
368 detection provided by the system (as confirmed by 120 consecutive correct detections); and the time to take
369 action, t_2 , between the system fault detection and the fault consolidation (beginning of steady-state phase).
370 According to the definitions of t_1 and t_2 , the value of $(t_1 + t_2)$ is constant and represents the duration of the

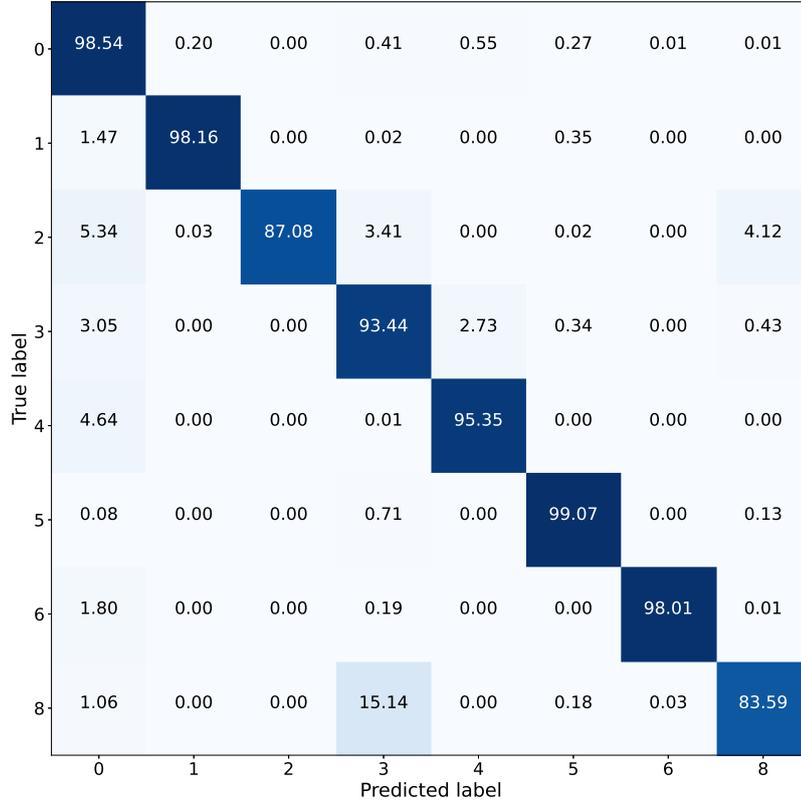


Figure 7: Confusion matrix showing final test results (in percentages) for Experiment 3.

Class	Initial Normal	Transient	Steady-state	All
0	127 (129)	-	-	-
1	37 (38)	36 (38)	30 (38)	36 (38)
2	10 (12)	8 (12)	1 (9)	7 (12)
3	-	-	23 (29)	-
4	-	-	74 (84)	-
5	1 (4)	99 (103)	85 (101)	0 (4)
6	45 (47)	44 (47)	36 (47)	43 (47)
8	25 (27)	24 (27)	17 (20)	24 (27)

Table 9: Event-level successful detection/classification results (and total number of events between parenthesis) for each event phase and fault class. Column “All” includes the events where the model detected the fault during its transient phase or during the steady-state period, and correctly identified the initial normal phase preceding the fault. Blank entries indicate lack of particular data in 3W dataset. The classifier used was the one of the model 3.3.

371 transient phase. So, as t_1 decreases (which indicates less time to detect faults), t_2 decreases.

372 The average values of t_1 and t_2 for each fault class in the test set are shown in Table 10, as well as the
373 average percentage values of these parameters with respect to the total transient-phase duration. Classes 0,
374 3, and 4 are omitted here as their transient phase is not specified in the 3W database. These results show
375 that the system can not only act as a fault classifier, but it can also anticipate the failure during its early
376 stage. In the worst scenario, the system detects the failure within 12% of its transient duration, giving an

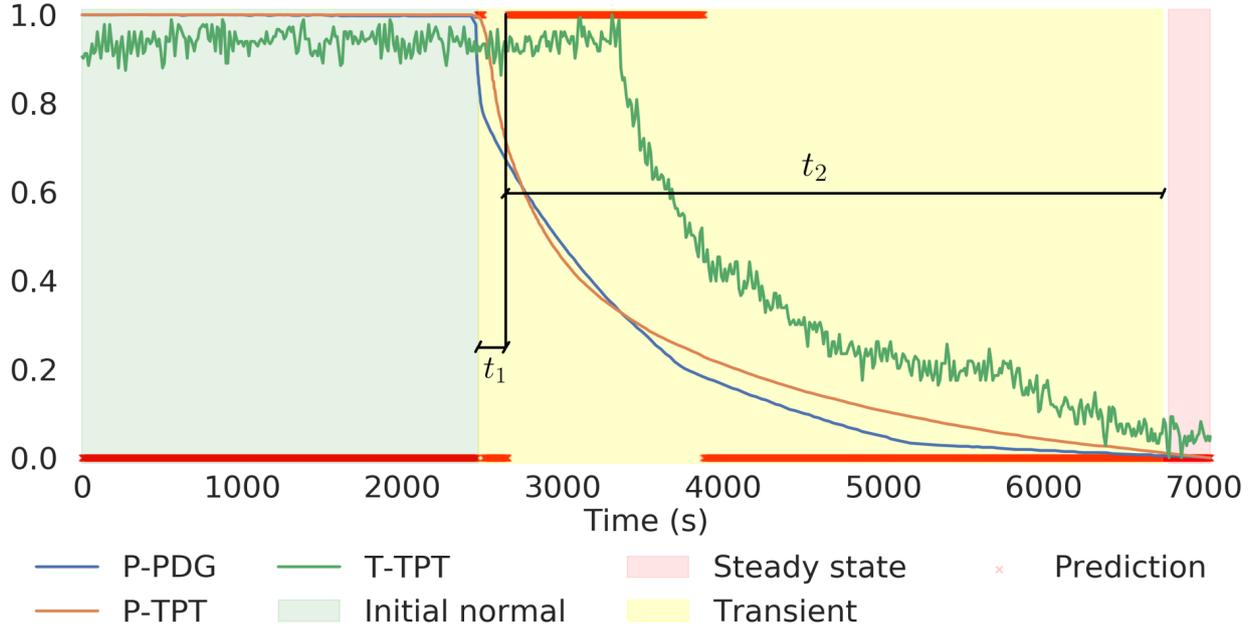


Figure 8: Example of a real instance of Class 2 alongside the system classification (red crosses) for each window sample: ‘0’ and ‘1’ outputs correspond to normal and faulty states, respectively.

377 additional time for the operators to intervene and prevent major production losses..

Class	t_1 [s]	t_2 [s]
1	293 (2.76%)	39804 (97.24%)
2	267 (8.77%)	4063 (91.23%)
5	6 (0.19%)	4966 (99.81%)
6	27 (3.69%)	6895 (96.31%)
8	2865 (11.09%)	15742 (88.91%)

Table 10: Average delay (t_1) and anticipation (t_2) intervals, in seconds, for the proposed multiclass Model 3.3 system. The number designated in parenthesis is the percentage of the corresponding time interval with respect to the total transient-phase duration. Classes 0, 3, and 4 are omitted here as their transient phase is not specified in the 3W database.

378 7. Conclusion

379 This paper described a full methodology for detecting and classifying faulty events during the practical
380 operation of O&G production wells and lines. Seven fault types were considered along with the normal
381 operation state. The developed system uses a classifier based on the random forest algorithm and a Bayesian
382 non-convex optimization strategy to tune the system hyperparameters.

383 Three experiments were devised to evaluate system capability and robustness in different fault detec-
384 tion/classification scenarios: Experiments 1 and 2 consider the binary normal \times faulty conditions, where the
385 faults are treated altogether and individually, respectively; Experiment 3 addresses the multiclass scenario,
386 where the system performs simultaneous fault detection and classification, which is best for practical usage.

387 In the multiclass configuration, for instance, overall accuracy results above 94% indicate successful perfor-
388 mance of the proposed system in detecting and classifying all faults types, thus reducing risk and production
389 losses in a real operation scenario. Alongside the high accuracy, the system also achieved a short detec-
390 tion delay, identifying the fault before completing 88% (in average) of its transient period, thus providing
391 additional time to the operator to mitigate associated damages.

392 Acknowledgments

393 This work was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior
394 - Brasil (CAPES) - Finance Code 001, CNPq, and Petrobras. The authors would also like to thank the
395 Petrobras team from the *Unidade de Negócios-Espírito Santo* for making the 3W dataset publicly available.

396 References

- 397 [1] T. Segaran and J. Hammerbacher, *Beautiful Data: The Stories Behind Elegant Data Solutions*, O'Reilly Media, 2009.
398 [2] M. Hilbert and P. López, "The world's technological capacity to store, communicate, and compute information," *Science*,
399 vol. 332, pp. 60–65, 2011.
400 [3] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer-Verlag, New York, 2006.
401 [4] S. Theodoridis and K. Koutroubas, *Pattern Recognition*, Academic Press, 4th ed., 2009.
402 [5] Y. S. Abu-Mostafa, M. Magdon-Ismael, and H.-T. Lin, *Learning from Data AML Book*, 2012.
403 [6] A. K. Jardine, D. Lin, and D. Banjevic, "A review on machinery diagnostics and prognostics implementing condition-based
404 maintenance," *Mechanical Systems and Signal Processing*, vol. 20, no. 7, pp. 1483–1510, Oct. 2006.
405 [7] I. H. Witten, E. Frank, and M. A. Hall, *Data Mining: Practical Machine Learning Tools and Techniques*, Morgan Kauf-
406 mann, 4th ed., San Francisco, 2016.
407 [8] W. Liu, B. Tang, J. Han et al., "The structure healthy condition monitoring and fault diagnosis methods in wind turbines:
408 A review," *Renewable and Sustainable Energy Reviews*, vol. 44, pp. 466–472, Apr. 2015.
409 [9] A. Grall, C. Bérenguer, and L. Dieulle, "A condition-based maintenance policy for stochastically deteriorating systems,"
410 *Reliability Engineering and System Safety*, vol. 76, no. 2, pp. 167–180, May 2002.
411 [10] R. C. M. Yam, P. Tse, L. Li et al. "Intelligent predictive decision support system for condition-based maintenance,"
412 *International Journal of Advanced Manufacturing Technology*, vol. 17, no. 5, pp. 383–391, Feb. 2001.
413 [11] G. M. Xavier and J. M. Seixas, "Fault detection and diagnosis in a chemical process using long short-term memory recurrent
414 neural network," *Proc. International Joint Conference on Neural Networks*, pp. 1–8, Rio de Janeiro, Brazil, 2018.
415 [12] A. Widodo and B.-S. Yang, "Support vector machine in machine condition monitoring and fault diagnosis," *Mechanical
416 Systems and Signal Processing*, vol. 21, no. 6, pp. 2560–2574, Aug. 2007.
417 [13] T. Helmy, A. Fatai, and K. Faisal, "Hybrid computational models for the characterization of oil and gas reservoirs," *Expert
418 Systems with Applications*, vol. 37, no. 7, pp. 5353–5363, July 2010.
419 [14] R. E. V. Vargas, C. J. Munaro, P. M. Ciarelli, A. G. Medeiros, B. G. Amaral, D. C. Barrionuevo, J. C. D. Araújo, J.
420 L. Ribeiro, L. P. Magalhães, "A realistic and public dataset with rare undesirable real events in oil wells," *Journal of
421 Petroleum Science and Engineering*, vol. 180, pp. 62–77, Oct. 2019.
422 [15] J. Snoek, H. Larochelle, and R. P. Adams, "Practical Bayesian Optimization of Machine Learning Algorithms", *Proc.
423 International Conference on Neural Information Processing Systems*, pp. 2951–2959, Nevada, 2012.
424 [16] I. H. Santos, A. R. M. Gaban, A. A. Lima, T. de M. Prego, B. D. Barros, M. A. Marins, E. A. B. da Silva, M. L. R.
425 Campos, and S. L. Netto, "On the automatic identification of critical production systems from a production loss database,"
426 *Proc. Rio Oil & Gas Expo and Conference*, pp. 1–7, Rio de Janeiro, 2018.
427 [17] I. H. Santos, H. F. Lisboa, T. de S. Feital, M. M. Câmara, R. M. Soares, M. A. Marins, B. D. Barros, T. de M. Prego,
428 A. A. Lima, and S. L. Netto, "Hydrate failure detection in production and injection lines using model and data-driven
429 approaches," *Proc. Rio Oil & Gas Expo and Conference*, pp. 1–13, Rio de Janeiro, 2018.
430 [18] W. McKinney, "Data Structures for Statistical Computing in Python," *Proceedings of the 9th Python in Science Confer-
431 ence*, pp. 56–61, Texas, Jun. 2010.
432 [19] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V.
433 Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine
434 Learning in Python," *Journal of Machine Learning Research* vol. 12, pp. 2825–2830, Oct. 2011.
435 [20] S. Kokoska and D. Zwillinger, *CRC Standard Probability and Statistics Tables and Formulae*, CRC Press, Boca Raton,
436 1999.
437 [21] L. T. Decarlo, "On the meaning and use of kurtosis," *Psychological Methods*, vol. 2, no. 3, pp. 292–307, 1997.
438 [22] D. C. Hoaglin and F. Mosteller, *Understanding Robust and Exploratory Data Analysis*, Wiley-Interscience, New York,
439 2000.

- 440 [23] C. C. Aggarwal, A. Hinneburg, and D. A. Keim, "On the surprising behavior of distance metrics in high dimensional
441 spaces," *Proc. International Conference on Database Theory*, pp. 420–434, Berlin, Jan. 2001.
- 442 [24] C. Feng, H. Wang, N. Lu, T. Chen, H. He, Y. Li and X. M. Tu, "Log-transformation and its implications for data analysis,"
443 *Shanghai Archives of Psychiatry*, vol. 26, no. 2, pp. 105–109, 2014.
- 444 [25] A. Ambrus, N. Saadallah, S. Alyaev, and F. Iversen, "Automatic detection of anomalous drilling operations using machine
445 learning methods and drilling process simulations," *Oil Gas European Magazine*, vol. 45, pp. 15–16, 2019.
- 446 [26] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, Oct. 2001.
- 447 [27] R. Díaz-Uriarte and S. Alvarez de Andrés, "Gene selection and classification of microarray data using random forest,"
448 *BMC Bioinformatics*, vol. 7, no. 1, pp. 3–16, Jan. 2006.
- 449 [28] M. Belgiu and L. Drăgut, "Random forest in remote sensing: A review of applications and future directions," *ISPRS*
450 *Journal of Photogrammetry and Remote Sensing*, vol. 114, no. 1, pp. 24–31, Apr. 2016.
- 451 [29] J. Jia, Z. Liu, X. Xiao, B. Liu and K. Chou, "pSuc-Lys: Predict lysine succinylation sites in proteins with PseAAC and
452 ensemble random forest approach," *Journal of Theoretical Biology*, vol. 394, no. 1, pp. 223–230, Apr. 2016.
- 453 [30] OSI Soft, *PI System*, [Online]. Available at: <https://www.osisoft.com/pi-system/>. [Access in 08/13/2019].
- 454 [31] Schlumberger, *OLGA Dynamic Multiphase Flow Simulator*, [Online]. Available at:
455 <https://www.software.slb.com/products/olga>. [Accessed in 08/13/2020].
- 456 [32] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, Springer-Verlag, New York, 2001.