



**HAL**  
open science

## Classification of voltage sags causes in industrial power networks using multivariate time-series

Maria Veizaga, Claude Delpha, Demba Diallo, Sophie Bercu, Ludovic Bertin

### ► To cite this version:

Maria Veizaga, Claude Delpha, Demba Diallo, Sophie Bercu, Ludovic Bertin. Classification of voltage sags causes in industrial power networks using multivariate time-series. IET Generation, Transmission and Distribution, 2023, 10.1049/gtd2.12765 . hal-04012128

**HAL Id: hal-04012128**


**<https://hal.science/hal-04012128>**

Submitted on 2 Mar 2023

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

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

# Classification of voltage sags causes in industrial power networks using multivariate time-series

Maria Veizaga<sup>1,3</sup> | Claude Delpha<sup>1</sup>  | Demba Diallo<sup>2</sup> | Sophie Bercu<sup>3</sup> | Ludovic Bertin<sup>3</sup>

<sup>1</sup>Université Paris Saclay, CNRS, CentraleSupélec, L2S, Gif Sur Yvette, France

<sup>2</sup>Université Paris Saclay, CNRS, CentraleSupélec, GeePs, Gif Sur Yvette, France

<sup>3</sup>EDF Lab Saclay, Palaiseau, France

## Correspondence

Claude Delpha, Université Paris Saclay, CNRS, CentraleSupélec, L2S, Gif Sur Yvette, France.  
Email: [claudedelpha@universite-paris-saclay.fr](mailto:claudedelpha@universite-paris-saclay.fr)

## Funding information

Association Nationale de la Recherche et de la Technologie, Grant/Award Number: 2019-12090

## Abstract

Voltage sags are the most frequent and impactful disturbances in industrial power grids, leading to high financial losses for industrial clients. The identification of the cause and its relative location is crucial for the contractual relation between the energy provider and the industrial customers. This paper proposes a methodology to identify the origins of voltage sags based on instantaneous symmetrical components and dynamic time warping. Short-Time Fourier and Fortescue transform are implemented in the pre-processing step using the voltage and current waveforms. Then, a distance-based classification strategy to identify the sources of voltage sags is used. It relies on a four-dimension time-series signature used as features. Moreover, a confidence index associated with the classification output is provided. The proposal offers an easy implementation in industrial applications with no previous recorded data. It has the benefit of using a reduced-size reference database entirely composed of synthetic data. The main advantages of the proposed method are its generalization capabilities and the possibility of raising an alert based on the confidence index. The obtained classification accuracy on synthetic data with seven causes is 100%. The method reaches a classification F1-score higher than 99% with field measurements representing five classes obtained from three different industrial sites.

## 1 | INTRODUCTION

The demand for power quality analysis in industrial networks has increased over the past decades. One of the main reasons is the financial impact of poor power quality. Among the electrical disturbances that can affect industrial networks, voltage sags are the most frequent and severe [1–3]. Understanding the origins of voltage sags is essential in the diagnosis process to implement corrective solutions and avoid financial losses due to production line downtime, damaged equipment, and wasted product. The three leading causes of voltage sags are: line faults, induction motor starting, and transformer energizing [4, 5], which can occur at the distribution network level or at the industrial site itself. Identifying the cause and its relative location is crucial for the contractual relationship between the energy provider and the industrial customers. Finding the voltage sag source allows the implementation of cost-effective and well-adapted

solutions to eliminate the sources of voltage sags (corrective solutions) or to protect the most sensitive equipment (mitigation solutions). For instance, if the sags are due to events inside the industrial site, corrective solutions can be more easily implemented according to the event causing the sag. Similarly, if the sags are due to events in the distribution system outside the maximum contract limits, the industrial customer may request financial compensation and compliance with the contract tolerances from the power utility operator. However, if the events are due to events in the distribution network within the contract limits, the industrial customer will be advised to implement targeted solutions on his site to protect only the most sensitive equipment and avoid financial losses.

Experts in power quality can analyse the data recorded by specialised monitoring devices installed on-premise during a limited amount of time, and provide a diagnosis. However, access to such expertise is not always available, and when it is,

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. *IET Generation, Transmission & Distribution* published by John Wiley & Sons Ltd on behalf of The Institution of Engineering and Technology.

it can be time-consuming and expensive. Therefore, the development of a methodology to automatically identify sources of voltage sag is a critical issue in the power quality diagnosis of industrial power systems. The classification methodology proposed in this work aims to identify the causes of voltage dips. This problem is different from the classification of voltage dip types, also called voltage dip characterisation, which has been abundantly addressed in the literature [6–9]. Indeed, a classification approach based on the causes of voltage sags may be more difficult but generally more relevant for the diagnosis of the electrical system, as pointed out by Bollen et al. in [10].

This paper proposes a method to identify the causes of voltage sags. It is based on a four-dimension time series signature. We consider Short-Time Fourier Transform (STFT), and instantaneous symmetrical components (Fortescue transform) to obtain distinctive signatures and use a dependent Dynamic Time Warping distance-based classification strategy. Our goal is to identify the cause and the relative location of a given voltage sag. The latter is equally important to allow industrial customers to determine whether the voltage drop is due to an event in the distribution grid or the industrial network, and therefore, they would be able to implement an appropriate countermeasure. Thus, we define the following seven causes of voltage dips affecting industrial sites: (a) upstream balanced faults, (b) upstream unbalanced faults, (c) downstream balanced faults, (d) downstream unbalanced faults, (e) upstream transformer energizing, (f) downstream transformer energizing, (g) downstream motor starting. In addition, the algorithm also provides a confidence index associated with the classification output, thus improving its reliability.

This paper is organized as follows. Section 2 describes the problem statement and related works, and Section 3 presents the characteristics of the main voltage sag causes and the model of the industrial network used for data generation. Section 4 details the classification method. Section 5 is dedicated to the performance analysis of the proposal in terms of class separability, robustness to noise, and fundamental frequency variations. Finally, in Section 6, a minimum database size is obtained, the results on real data are presented, and the classification accuracy analysis is detailed. The main conclusions are given in Section 7.

## 2 | PROBLEM STATEMENT AND RELATED WORKS

Most of the existing analysis methods in this field are mainly based on scalar feature extraction strategies, which are used as inputs to the classifier. Examples of such work include the following. [11, 12] proposed a solution based on S-transform (ST) in combination with Extreme Learning Machine (ELM) as a classifier. In [13], the authors used the Empirical Mode Decomposition (EMD) and the Hilbert transform (HT) combined with a Probabilistic Neural Network (PNN). In [14–17], the authors implemented variants of the Wavelet transform (WT) for feature extraction and different classifiers or combinations of classifiers such as Support Vector Machine (SVM), Random Forest (RF), Multilayer Neural Network (MLP), PNN and Naive Bayes. In

[18], Nagata et al. presented a method for the detection and classification based on High Order Statistics (HOS) and Artificial Neural Networks (ANN).

One of the main limitations of these approaches is their dependence on the choice of features, which should be done carefully to avoid the loss of relevant information in the process. In order to extract robust features, it is necessary to have a large and representative dataset for training, which can be a particularly limiting aspect, as voltage dips are disturbances of limited occurrence and therefore challenging to measure. Thus, building up an extensive training data set is a major obstacle to implementing these methods in real industrial applications. In addition, the transformations used in the pre-processing step and the retained features rarely take into account the physical properties of the events. Thus, the generalisation capabilities of these algorithms are compromised, as it is difficult to provide guarantees on the behaviour of the selected features when applied to new data from the training data set. This is a well-known problem in machine learning, known as Domain Adaptation [19], defined as the ability to successfully apply an algorithm trained in a source domain to a different but related target domain. In our research area, this problem was highlighted by Bollen et al. in [10] when a model trained with a synthetic dataset was tested on measured data, and provided non-acceptable results.

Classical machine learning algorithms only deal with scalar features in the examples mentioned above. This approach may be well suited to the analysis of steady-state disturbances. However, the time dependence of electrical waveforms is important for analysing short-duration disturbances such as voltage dips. This is because the information related to the underlying causes of voltage dips is encoded through the entire duration of the event, and extracting scalar features involves a risk of information loss [20]. For this purpose, a classification approach based on time series seems more relevant for a more efficient analysis of voltage sags. Time series classification is an area of machine learning that has developed in recent years, and new algorithms have recently appeared [21, 22]. Recent approaches to identifying the causes of voltage drops take advantage of temporal waveforms by applying deep learning algorithms. In [23] a self-supervised voltage sag source identification method based on Convolutional Neural Networks (CNN) and an autoencoder has been proposed. A bidirectional long short-term memory network (Bi-LSTM) was proposed in [24] and [25]. Indeed, these models have no feature extraction step as deep neural networks are designed to extract their own features during the training stage. Even if this can be an advantage when access to expert knowledge is not available, these approaches are highly data-driven, and their performance strongly depends on the characteristics of the training dataset (size, diversity etc.).

Although deep learning approaches are popular, there are other time series classification methods that are less dependent on the size of the training dataset. For instance, 1-Nearest Neighbor classifier with Dynamic Time Warping (1NN-DTW) is a recommended benchmark due to its simplicity and hard-to-beat accuracy [22]. However, one of the disadvantages of this method is the small number of neighbors involved in the

classification, which may decrease its robustness. The time alignment difference between compared signals is handled by the DTW algorithm [26]. Youssef et al. [27] proposed a power quality event classification method using the DTW algorithm, Walsh, and Fast Fourier transform. In [28], we presented a voltage sag cause classification method based on symmetrical components and 1NN-DTW as a classifier.

It is also important to notice that compared to the majority of data-driven approaches in the literature, few proposals include expert knowledge in their solutions by extracting meaningful features from the electrical point of view [29, 30]. Usually, these approaches also implement rule-based strategies for the classification stage by setting thresholds, which can prevent algorithms from generalizing optimally when applied to different grid configurations. Another way of incorporating electric-based concepts into diagnosis algorithms is to use some of the tools usually applied in this field. For instance, the Clarke transform has been used in [8, 31] for the characterization of voltage sags and swells and in [32] for fault classification and voltage sag parameter computation. The Fortescue transform is also widely used to analyze unbalanced power systems. It has been used for the characterization of voltage sags [7, 33, 34], and for determining the relative location of voltage sag origins [35, 36]. But, to the best of our knowledge, it has not been applied for the classification of the causes of voltage sags.

### 3 | INDUSTRIAL NETWORK MODEL AND MAIN VOLTAGE SAG SOURCES

#### 3.1 | Main voltage sag sources

The leading causes of voltage sag reported in the literature are: line faults, energizing of transformers, and starting of induction motors, as they considerably increase the amplitude of the absorbed current. Moreover, current and voltage waveforms depend on the event causing the voltage sag and on the loads in the industrial network. It should also be noted that the characteristics of the distribution network, in particular the power available at the point of common coupling (PCC), has a significant impact on the magnitude of the voltage drop.

##### 3.1.1 | Line faults

The most common cause of voltage sags is line faults, which can be originated either from the upstream or downstream networks. They can be caused by lightning or any object in contact with energized lines. There are different types of faults depending on the affected phases. Faults can also be classified in symmetrical and asymmetrical. Three-phase faults (LLL) or three phase-to-ground faults (LLLG) are considered symmetrical, while single phase-to-ground faults (L-G), double-phase faults (LL), and double phase-to-ground faults (LL-G) are considered asymmetrical. During a fault occurrence, the voltage drops because of the high inrush current. The voltage is restored

in the healthy part of the grid as soon as the faulted feeder is detected and de-energized by the corresponding protective device (upstream case) or as soon as a protective device (e.g. circuit breaker or fuse) detects the over-current and disconnects the faulted asset (downstream case).

##### 3.1.2 | Transformer energizing

In transmission and distribution grids, the energizing of a transformer usually follows the end of a protection cycle: when the faulted supply is disconnected, the loaded transformers are energized. Transformer energizing can also occur in industrial networks downstream of the monitoring point. The magnetic flux may exceed the saturation limit when energizing the transformer, creating a high inrush current and a voltage drop. Core saturation affects the frequency spectrum of waveforms with additional harmonics during the transient. A significant level of voltage unbalance also affects all three phases. Another characteristic is the voltage exponential recovery, determined by the time constant with which the residual flux decreases in the transformer core [37].

##### 3.1.3 | Induction motor startup

Induction motor online startup creates a high inrush current that can reach 7–10 times its nominal value, causing a measurable voltage sag at the monitoring point. As voltage drops, current rapidly increases. Voltage recovery is achieved at the end of the startup as the current slowly decreases. The total startup time depends on motor characteristics (power, inertia, torque, load type, ...), ranging from less than one second for small motors to a few seconds for large motors.

#### 3.2 | Industrial network simulation model description

The simulated case study is designed on the EMTP® software [38] (Electromagnetic Transients Program), which is a simulation tool for load-flow, steady-state, and time-domain analysis for power systems. The electric model is composed of two sub-networks: the industrial network and the distribution network, as displayed in Figure 1. Note that the monitoring device is placed downstream of the main MV/LV transformer. The location of the events is relative to the monitoring point. Thus, events occurring in the distribution network are qualified as upstream, and those occurring in the industrial network as downstream.

##### 3.2.1 | Industrial network

For our study, the nominal voltage of the industrial network is 400 V, and the frequency is 50 Hz. The site is connected to a  $Dy11n$  21 kV/400 V transformer of 400 kVA. In order to model

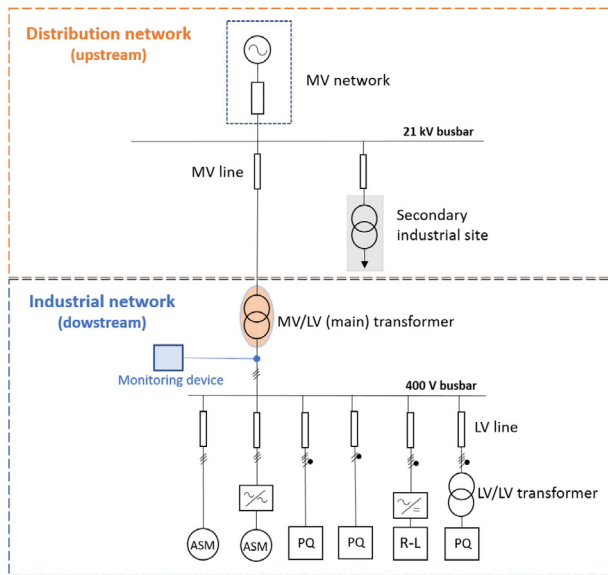


FIGURE 1 Simplified diagram of the industrial network model

the behaviour of a reduced but representative LV industrial site, the following loads were included:

- *Induction motors* with rated power varying from 22 to 110 kW, with variable inertia and torque.
- *Variable speed drivers* with scalar control.
- *Three-phase rectifiers* designed as 6-pulse rectifiers feeding 30 kW DC loads.
- *Three-phase loads* of 30 kW and  $\cos\phi$  of 0.9.
- *Single-phase loads* of 5.5 kW and  $\cos\phi$  of 0.6.
- *Three-phase isolation transformers* (400 V/400 V) with rated power varying from 100 to 250 kVA with  $Dyn$  and  $YNyn$  winding connections.

The loads are not all simultaneously connected. For instance, the large motors and transformers are connected only for the motor startup and the transformer energizing scenarios. The nominal power consumption of the modeled industrial site is around 350 kW. In order to emulate the critical state where the network is susceptible to experiencing voltage sags, the main 400 kVA transformer is not optimally designed. Some of the devices are sources of permanent disturbances frequently present in industrial networks:

- *Harmonics.* Variable speed drivers and 6-pulse rectifiers generate current harmonics, mostly of orders 5 and 7.
- *Unbalance.* The presence of single-phase loads simulates the injection of a certain level of unbalance.

### 3.2.2 | Distribution network

In the distribution network (upstream), the voltage level is 21 kV. The nominal power of the main busbar varies in the range [10 to 100] MVA, and the ratio of the reactance to the resistance

of the system is 0.1. The length of the MV lines ranges between 0 and 30 km. Two feeders connect the 21 kV busbar to the main industrial network and to a secondary site. The latter is composed of  $Dyn11$  transformers whose rated power varies between 500 and 1250 kVA. The transformers are connected in parallel, feeding loads of 300 kW.

## 3.3 | Voltage sag generation

### 3.3.1 | Line faults

Different line fault types are generated: single line-to-ground (L-G), line-to-line (LL), double line-to-ground (LL-G), and three-phase (LLL) faults. The magnitude of voltage sag mainly depends on the fault type, the distance to the fault (line length), and the value of the ground fault resistance. The duration of the sag in this case depends on the action of the protection equipment. Typical fault clearance time varies between 3 and 30 cycles (60–600 ms in a 50 Hz grid)[4]. This type of event is generated at three locations: the secondary site's feeder MV line (upstream), the 400 V busbar of the main industrial network (downstream), and the secondary side of an isolation transformer (downstream).

### 3.3.2 | Transformer energizing

The magnitude of the voltage sag caused by this event depends on the transformer's power and its core flux initial state (different from zero when the core has not been entirely demagnetized before re-energizing). The duration of the sag depends on the transformer's characteristics. This event is generated at two locations: the secondary side (upstream) and at the 400 V busbar of the industrial network (downstream).

### 3.3.3 | Direct online motor startup

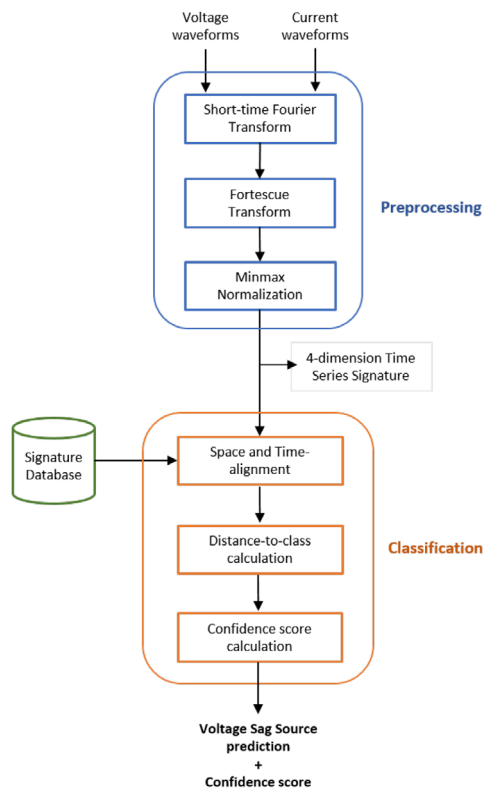
The magnitude of the voltage sag caused by a motor startup depends on its power, torque, and total inertia. These parameters also determine the entire duration of the event and the sag as a consequence. It is generated by directly connecting the induction motor to the 400 V busbar of the industrial network (downstream).

The generated voltage sags vary between 10% and 98% (residual voltage). Although a voltage sag is defined as the reduction of voltage RMS values lower than 90%, shallow voltage drops are also included in the dataset. Their identification can be more difficult than deep voltage sags since the variations on voltage and current can be very low.

## 4 | CLASSIFICATION OF VOLTAGE SAG CAUSES

The proposed classification method described in Figure 2 mainly relies on the pre-processing step. The voltage and





**FIGURE 2** Flowchart of the proposed voltage sag source classification method

current waveforms, measured at the monitoring point installed at the secondary of the main transformer, are used as inputs. The first transformation in the pre-processing step is the Short-Time Fourier Transform (STFT), which accomplishes two tasks: (a) decomposition of the signals into their harmonic content and (b) computation of the complex values of real waveforms. The second task is essential as the Fortescue transform is defined in the complex domain. Then, the Fortescue transform converts the three instantaneous phasors into three instantaneous symmetrical components (positive, negative, and zero-sequence). Four instantaneous symmetrical components are retained: (a) positive-sequence voltage (harmonic 1), (b) negative-sequence voltage (harmonic 1), (c) positive-sequence voltage (harmonic 2), and (d) positive-sequence current (harmonic 1). Finally, these four components are normalized using min-max normalization.

The set of four-time series represents the voltage sag signature (four-dimensional signature). The signature of each of the seven defined classes can be visually identified and distinguished from the others.

The pre-processing stage is based on the one briefly presented in [28]. The main contributions on a new optimized and more robust distance-based classification strategy. For this purpose, a database containing the signatures of labelled events belonging to each class was constructed. The classification of a new event is based on the analysis of the distances between its signature and those in the database. A spatio-temporal alignment step is first performed between the new signature and

those in the database. Then, the distance of the new signature to each class in the database can be calculated. The predicted label assigned to the event will correspond to the closest class. Finally, two confidence scores on the classification output are calculated and provided.

## 4.1 | Pre-processing

### 4.1.1 | Short-time fourier transform

STFT [39] is a technique used for the analysis of the frequency content variation of a non-stationary signal. It is obtained by applying the Discrete Fourier transform (DFT) over the signal through a sliding window of length  $W_L$ . The window overlap between signal segments compensates for the signal attenuation at the window edges. The DFT of each segment is saved into a matrix containing the magnitude and angle for each data point in time and frequency.

The matrix is defined as:

$$STFT(x) = [X_0(f), X_1(f), \dots, X_T(f)], \quad (1)$$

with,

$$X_m(f) = \sum_{n=-\infty}^{\infty} x(n)g(n-sR)e^{-j2\pi fn}, \quad (2)$$

where  $n$  is the length of the discrete signal  $x(n)$ ,  $X(f)$  is the DFT of the windowed signal centered at time  $sR$ ,  $R$  is the distance between two adjacent windows and  $g(n)$  is the window function.

The parameters are set as:  $W_L = F_s/F$ , with  $F_s$  the sampling frequency of  $x(n)$ ,  $F$  the nominal frequency (50 Hz),  $R = 1$  and  $g(n)$  is a rectangular window. The signal is decomposed into a set of frequency bands corresponding to its harmonics ( $f_1 = 50$  Hz,  $f_2 = 100$  Hz,  $f_3 = 150$  Hz...). The rectangular window avoids amplitude attenuation, and the minimum distance between adjacent windows (maximum overlap) gives the best possible time resolution.

The first reason we implement STFT to decompose the signal into its harmonic components is only to extract the harmonics of interest and avoid those affected by the industrial loads. For instance, harmonics of orders 5 and 7 are significantly affected by 6-pulse rectifiers (present in variable speed drivers and three-phase rectifiers). Their isolation increases the generalization capacity of the algorithm despite the different types of loads present in the industrial network. The fundamental frequency (harmonic of order 1) contains primary information for voltage sag analysis, as sags are low-frequency disturbances. The presence of even harmonics due to transformer energizing is also useful information for their identification. Therefore, the harmonic of order 2 is also extracted from the STFT matrix. The second reason for implementing STFT, is that the Fortescue transform to be applied next is defined in the complex domain. Real waveforms of voltage and current are converted into corresponding complex harmonic phasors.

## 4.1.2 | Fortescue transform

Fortescue transform [40] is a linear transformation used for the analysis of unbalanced three-phase power systems. It transforms an unbalanced set of three phasors ( $\underline{X}_A, \underline{X}_B, \underline{X}_C$ ) into a balanced set of three symmetrical components, that is, the positive ( $\underline{X}_+$ ), negative ( $\underline{X}_-$ ) and zero-sequence ( $\underline{X}_0$ ). The transform is defined as:

$$\begin{bmatrix} \underline{X}_0 \\ \underline{X}_+ \\ \underline{X}_- \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & \underline{\alpha} & \underline{\alpha}^2 \\ 1 & \underline{\alpha}^2 & \underline{\alpha} \end{bmatrix} \begin{bmatrix} \underline{X}_A \\ \underline{X}_B \\ \underline{X}_C \end{bmatrix}, \quad (3)$$

where  $\underline{\alpha} = e^{j\frac{2}{3}\pi}$  is the phasor rotation operator.

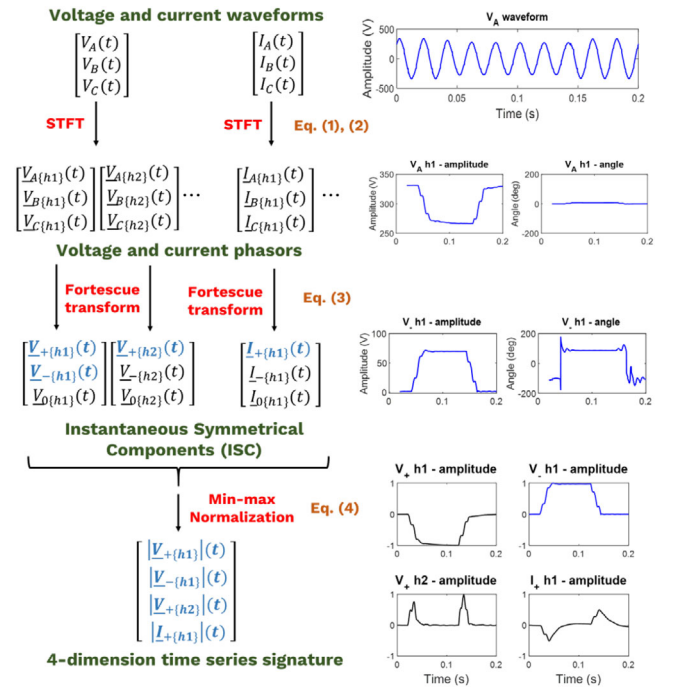
The Instantaneous Symmetrical Components (ISCs) of the harmonics extracted from voltage and current waveforms are calculated using the instantaneous complex values previously determined by STFT. Positive-sequence components represent the actual voltage and current being provided to the load. In a perfectly balanced system, negative and zero-sequence components are equal to zero. Zero-sequence is directly related to the grounding system and transformer winding connections. However, for upstream disturbances, the zero-sequence component is filtered [41] since MV/LV transformers of industrial sites in France usually have a  $Dy$  winding connection. Since the induction motors' windings are connected in  $\delta$  ( $D$ ) or ungrounded  $nye$  ( $Y$ ), no zero-sequence current is generated. Therefore, the zero-sequence voltage is not influenced by the induction motor either [33]. Finally, four ISCs are selected: (a) positive-sequence voltage (harmonic 1), (b) negative-sequence voltage (harmonic 1), (c) positive-sequence voltage (harmonic 2), and (d) positive-sequence current (harmonic 1). Their characteristics are described in more detail in Section 4.2.

## 4.1.3 | Minmax normalization

The selected ISCs constitute a four-dimension time series signature. We use the shape of each ISC as the main discriminant characteristic between classes (causes) of voltage sags. A min-max normalization is applied to each ISC to perform a shape-based time series classification. Minmax normalization with a re-scaling between  $[a, b]$  is defined in (4) with  $X$  being the ISC to be normalized,  $a = -0.5$  and  $b = 0.5$ .

$$X_{[a,b]} = \frac{X - \min(X)}{\max(X) - \min(X)} * (b - a) + a. \quad (4)$$

The ISCs are first re-scaled, then zero-centered by subtracting the value of the first point from the rest of the sequence. Each time series is defined between  $[-1, 1]$  at the end of this operation. Figure 3 illustrates the different steps of the pre-processing and the components obtained at each step.



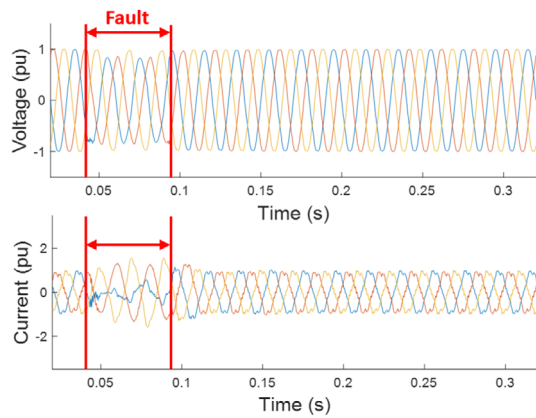
**FIGURE 3** Detailed steps of the pre-processing and obtained components for the calculation of the four-dimensional signatures

The four-dimension signatures of real records of voltage sags are illustrated: Figure 4 shows an upstream unbalanced fault, Figure 5 an upstream transformer energizing, and Figure 6a downstream motor startup.

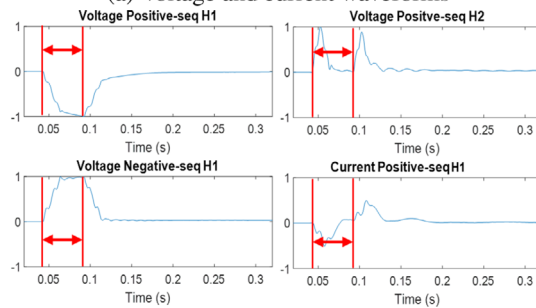
## 4.2 | Electrical interpretation of signatures

The four-dimensional signatures obtained at the end of the pre-processing stage can be interpreted from an electrical point of view. Each ISC reflects one or more electrical characteristics of the event.

1. *Voltage positive-sequence harmonic 1* component represents the voltage evolution during the sag, with very similar characteristics to the RMS three-phase voltage curves. The rapid drop is directly related to the onset of the voltage sag. For example, the quasi-square shape is characteristic of sags caused by faults, which have a rapid recovery after the fault is cleared. Sags caused by starting motors and energized transformers have, in contrast, a progressive recovery with characteristics similar to their RMS voltage curves.
2. *Voltage negative-sequence harmonic 1* component reflects the unbalanced nature of the voltage sag. For instance, the sudden increase and sustained relative high values during the sag are expected characteristics of an imbalanced line fault. On the contrary, a symmetrical fault will only present two peaks at the sag transients' beginning and end. A three-phase motor starting will show only one peak (with one or more lobes) at the first and only transient, as it causes a balanced

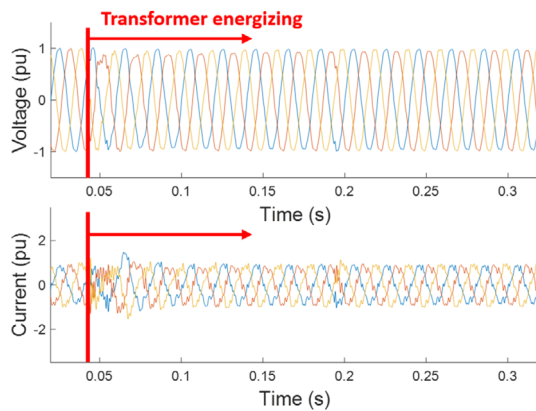


(a) Voltage and current waveforms

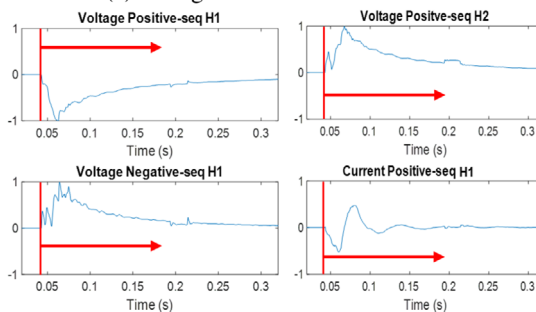


(b) Time series signature

**FIGURE 4** Voltage sag caused by an upstream non-symmetrical fault (real measurement data)

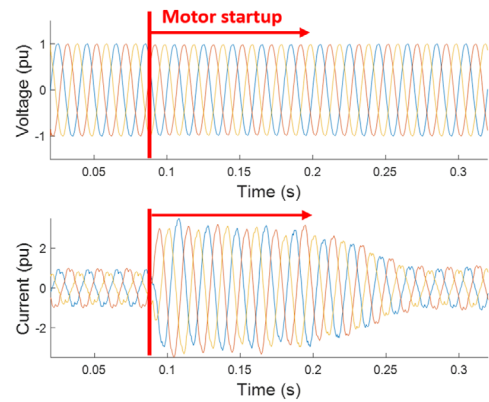


(a) Voltage and current waveforms

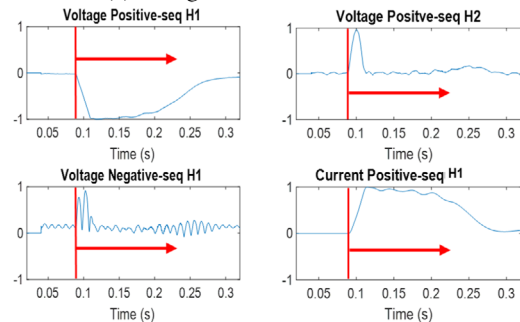


(b) Time series signature

**FIGURE 5** Voltage sag caused by an upstream transformer energizing (real measurement data)



(a) Voltage and current waveforms



(b) Time series signature

**FIGURE 6** Voltage sag caused by a downstream induction motor startup (real measurement data)

sag. Finally, a transformer energizing will at first cause a high increase and then a slow decrease until the voltage is recovered. This behavior is also expected as the transformers cause unbalanced voltage sags.

3. *Voltage positive-sequence harmonic 2* component is particularly useful to identify a voltage sag caused by a transformer energizing since it causes the onset of even harmonics. At the fault occurrence, two peaks are present at the beginning and the end of the sag. For motor starting, only one peak is visible at the beginning of the event, corresponding to a transient.
4. *Current positive-sequence harmonic 1* component is used to determine the relative location of the event (upstream or downstream of the monitoring device). For upstream events, it rapidly decreases at the event occurrence time. This phenomenon can be explained by the energy sink analogy presented in [42], where events such as faults or load connections consume high amounts of current and energy. Then, the current recovers and stabilizes during the voltage sag, this stage being visible during sags caused by faults. When the voltage is restored at the end of the sag, the current will reach a peak value higher than its nominal value. The magnitude of the peak depends on the duration and severity of the sag and the connected loads (motors power and inertia, DC bus capacity etc.). For downstream events, on the contrary, this ISC will increase and keep sustained values as long as the sag is on and will decrease proportionally to the voltage recovery curve.



### 4.3 | Signatures classification

The classification of the events is based on the similarity of the time series signatures. A similarity measure can be obtained by calculating the distance between two signatures. Classical distance measures such as Euclidean distance are not suitable when applied directly to similar time series that slightly differ in time or speed. Thus, an algorithm such as DTW [26], capable of handling time misalignment, is necessary. In the same way, spatial misalignment can be a problem when analyzing incomplete signals (due to the limitations of the monitoring device).

#### 4.3.1 | Space and time alignment

##### Space-alignment

The characteristics and parameters of the monitoring device are sometimes not well calibrated. This can cause incomplete event recordings, which can affect the signature-matching process. This can be especially problematic if the missing segment corresponds to the beginning of the event, as it is necessary to have at least one “healthy” period before the sag. This first period is used as a reference for zero-centering the signatures. An offset can be applied to correct space misalignment.

At first, we verify if there exists a risk of having less than one reference period due to an incomplete registered event by calculating the difference between the mean value of the first period and the last period for each ISC of the signature. If the difference  $\delta$  is higher than a minimum threshold of 0.05 pu, an offset can improve the alignment. A vector containing possible offsets is defined as  $X_{off} = -sign(\delta) * [0, 0.05, 0.1, \dots, \delta]$ , and the Euclidean distance is calculated between the ISC query plus the offset and the ISC reference. The minimal Euclidean distance gives the optimal offset, which is applied at the end of this stage.

##### Time-alignment

Dynamic Time Warping (DTW) [26] is a well-known algorithm used for handling time alignment differences between two univariate time series. The extension of this algorithm to multivariate time series can be achieved either by calculating the DTW distance for each dimension independently and adding the calculated distances (independent dynamic time warping or  $DTW_I$ ) or by calculating the DTW distance across all the dimensions simultaneously (dependent dynamic time warping or  $DTW_D$ ) [22]. In our case, the four dimensions of each signature are time-correlated, and a single optimal time alignment warping path is calculated for all the dimensions.

A local distance matrix of size  $L \times L$  is calculated,  $L$  being the length of the signature. Each element of the distance matrix is defined through a chosen distance measure, the Euclidean distance being usually the privileged choice. The equation calculating the elements  $e(i, j)$  of the distance matrix for two signatures  $q$  (query) and  $r$  (reference) of dimension  $G = 4$  with

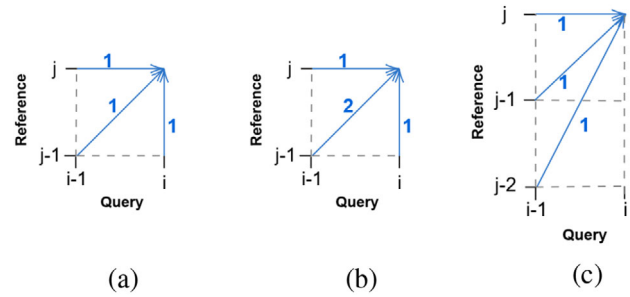


FIGURE 7 Local step patterns: (a) *symmetric1*, (b) *symmetric2*, (c) *asymmetric* [44]

indexes  $i, j \in [0, L]$  is given by

$$e(i, j) = \sum_{g=1}^{G=4} (q(i, g) - r(j, g))^2. \quad (5)$$

The optimal warping function  $\phi(n)$  defined in (6), is found through the minimization of the cumulative cost  $E$  defined in (7) obtained from the distance matrix. It is given by two integer vectors  $\phi_q(n)$ ,  $\phi_r(n)$  of same length  $N$  (with  $L \leq N \leq 2L$ ), mapping the time axis of the query  $q$  to the reference  $r$ . These vectors indicate the time alignment to be applied to both signatures  $q$  and  $r$  to all the dimensions.

$$\phi(k) = (\phi_q(n), \phi_r(n)) \quad \text{with } k = [0, N], \quad (6)$$

$$E_{min}(q, r) = \min_{\phi} \sum_{n=1}^N a(\phi_q(n), \phi_r(n)) w(n). \quad (7)$$

Finally, the normalized distance between the aligned query and reference  $D(q, r)$  is defined in (8). The factor  $1/N$  normalizes the total distance, regardless of the length  $N$ , which tends to increase when signatures are stretched or compressed in time.

$$D(q, r) = \frac{1}{N} \sum_{g=1}^{G=4} \sqrt{\sum_{n=1}^N (q(\phi_q(n), g) - r(\phi_r(n), g))^2}. \quad (8)$$

Local or global constraints can be applied to limit certain types of time series distortions, also called “singularities.” These constraints are included by modifying the weights  $w(n)$  used to calculate the cumulative cost.

Global constraints such as the Sakoe-Chiba band [26] or the Itakura parallelogram [43] limit the distance of the warping function to the main diagonal. However, these constraints are not adequate for our problem since they can prevent an optimal alignment of two similar events but with highly different duration, limiting time dilation or compression capabilities.

Local constraints include step patterns, which are more flexible but still can limit severe signal distortion. Figure 7 illustrates three well-known step patterns: *symmetric1*, *symmetric2* and *asymmetric* [44]. The best step pattern should allow minimal distance

**TABLE 1** Step pattern comparison

	<i>symmetric1</i>		<i>symmetric2</i>		<i>asymmetric</i>	
	$D$ (e-02)	$Z$ (%)	$D$ (e-02)	$Z$ (%)	$D$ (e-02)	$Z$ (%)
Intra-class	1.85	19.79	1.41	71.83	2.26	40.74
Inter-class	9.48	20.59	6.26	90.11	6.89	47.15

between signatures of the same class and maximal distance between signatures of different classes, with the distortion of the time series being minimal in both cases. The distance between two signatures is given by Equation (8). The distortion rate  $Z_{(\%)}$  is defined in Equation (9), with  $N$  being the warping path length and  $L$  the length of the original signature. Since the two compared signatures are equal in length, a perfect time alignment would correspond to a diagonal warping path of length  $N = L$  and  $Z_{(\%)} = 0$ . The maximum warping path length being  $N = 2L$ , the maximum value of  $Z_{(\%)} = 100$ .

$$Z_{(\%)} = \frac{N - L}{L}. \quad (9)$$

Table 1 shows the distance  $D$  and distortion rate  $Z$  for signatures belonging to the same class (intra-class) and between classes (inter-class) for the three-step patterns. The *symmetric1* step pattern has the best trade-off between minimal intra-class distance, maximal inter-class distance, and minimal distortion rate.

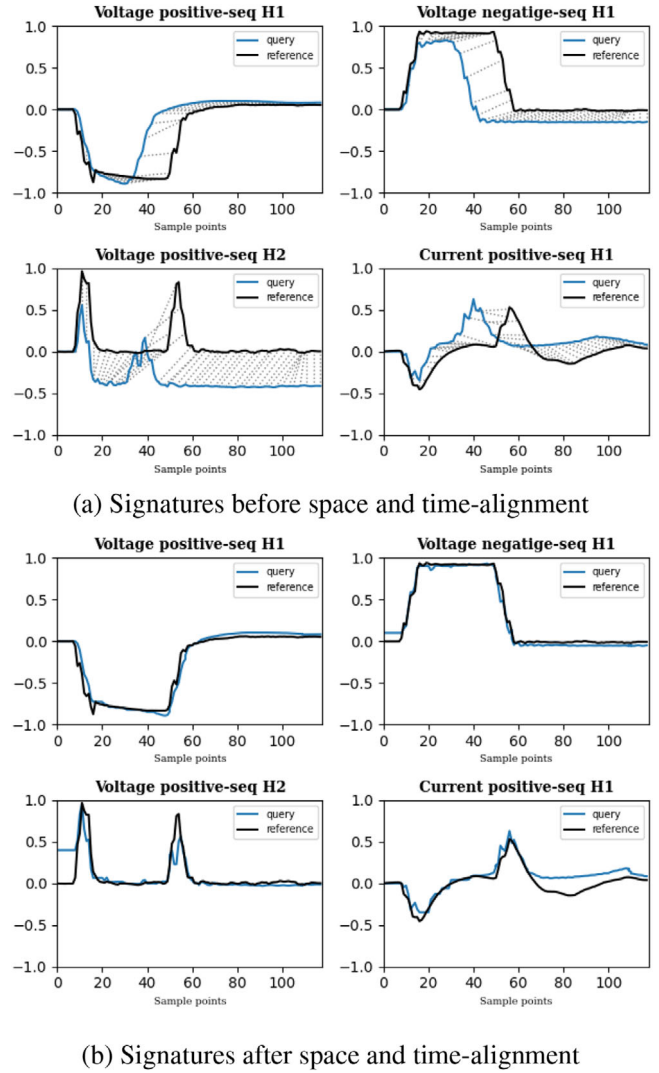
The *symmetric1* step pattern favors oblique steps over horizontal or vertical ones. This characteristic limits distortion compared to the *symmetric2* step pattern, which considers a vertical plus horizontal step equivalent to an oblique step. The *symmetric1* step pattern also allows a higher degree of dilation or compression compared to the *asymmetric* step pattern. The latter imposes a single match point for every point in a time series.

Figure 8 illustrates the result of the space and time-alignment step using the *symmetric1* step pattern. The query signature corresponds to a real voltage sag and the reference to a synthetic voltage sag.

### 4.3.2 | Distance-to-class calculation

The classification of a new voltage sag signature  $q^*$  is achieved by the mean distance to all the reference signatures  $r_k$  belonging to the class  $C_k$  of size  $M$ . In order to obtain a more robust estimator of the mean distance to each class  $d_k(q^*)$ , we define  $\hat{d}_k(q^*)$  using a bootstrapping approach, as described in (11). We define  $B$  as the total number of sub-samples  $X_{D,b}$ , extracted from the population of distances  $X_D = \{D(q^*, r_{k,m})\}$  with  $m = 1, 2, \dots, M$ , and  $\bar{d}_{k,b}$  as the mean of  $X_{D,b}$ .

$$\bar{d}_k(q^*) = \frac{1}{M} \sum_{m=1}^M D(q^*, r_{k,m}), \quad (10)$$



**FIGURE 8** Result of the space and time alignment step. The query corresponds to a real voltage sag and the reference to synthetic voltage sag. The space alignment corrects the error caused by the monitoring device

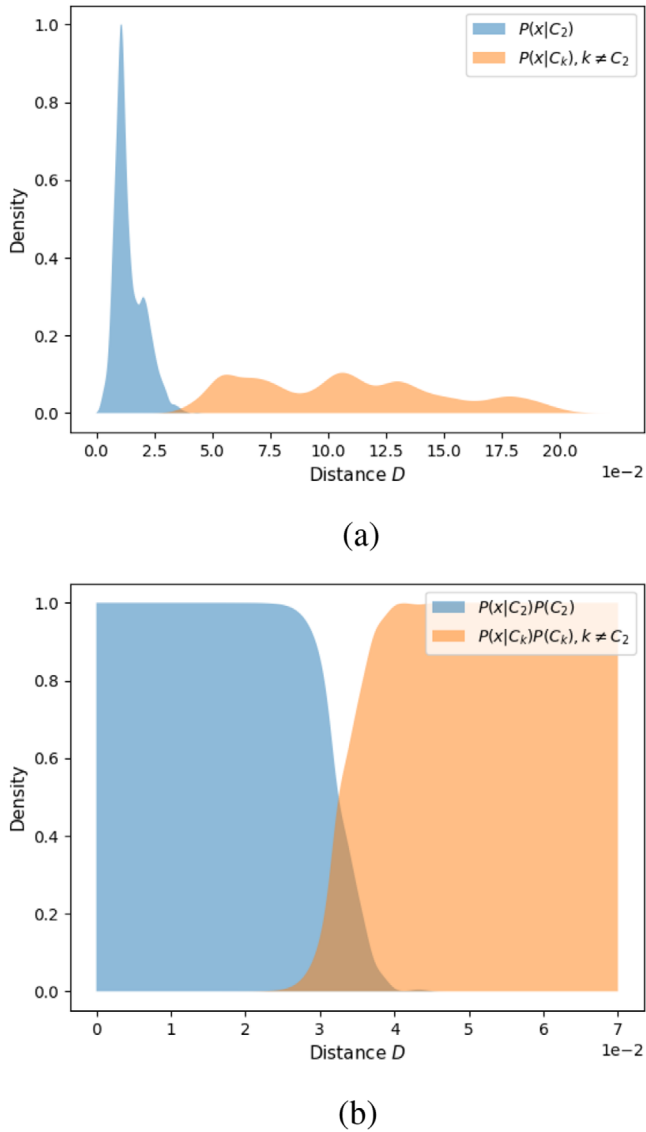
$$\hat{d}_k(q^*) = \frac{1}{B} \sum_{b=1}^B \bar{d}_{k,b}. \quad (11)$$

Once the distances  $d_k(q^*)$  between the new voltage sag event signature  $q^*$  to each class  $C_k$  are calculated, we obtain a vector containing the distances to the  $K$  classes. The closest class gives the label  $y^*$  assigned to the event such as:

$$y^* = \arg \min_{k \in \{1, \dots, K\}} ([d_1(q^*), d_2(q^*), d_k(q^*) \dots, d_K(q^*)]). \quad (12)$$

### 4.3.3 | Confidence score on the classification output

In order to validate the efficiency of the method, we propose to compute a confidence score associated with the classification output result. Therefore, we define two confidence



**FIGURE 9** Likelihood function (a) and output score (b) of a binary Naive Bayes classifier using KDE. New voltage sag signature  $q^*$  to the class  $C_2$ , given the distance to the class  $x = d_2(q^*)$

indexes: a probabilistic-based index (NB-KDE) and a relative distance-based index (RD).

#### Probabilistic (NB-KDE) index

The first confidence index is obtained through a set of  $K$  binary Naive Bayes classifiers using a Kernel Density Estimator (NB-KDE). Each binary classifier is trained using a one-vs-rest approach. According to Bayes theorem given in (13), where  $x = d_k(q^*)$ , the posterior probability  $P(y = C_k|x)$  is proportional to the likelihood function estimated by the KDE  $P(x|y = C_k)$  and the prior  $P(C_k)$ . Figure 9 shows the likelihood function and the output score of a binary Naive Bayes classifier using KDE.

$$P(y = C_k|x) \propto P(x|y = C_k)P(y = C_k). \quad (13)$$

The probabilistic index NB-KDE associated with the classification output given in (12) is the output provided by the binary classifier  $k = y^*$ , as defined in (14).

$$\text{NB-KDE} = P(x|y = C_{y^*})P(y = C_{y^*}). \quad (14)$$

The closer a new voltage sag signature  $q^*$  is to the reference signatures  $r_k$  of a given class  $C_k$  of the database, the higher the confidence index NB-KDE will be. As expected, the number of reference signatures in the database influences the estimation of the likelihood function. An extensive and variate database will provide a reasonable estimate of the likelihood function by the KDE, thus, a more accurate NB-KDE confidence index. On the contrary, a reduced database will produce a poor likelihood function estimation, and the NB-KDE confidence index will not be considered as reliable.

#### Relative distance-based (RD) index

The second confidence index is based on the calculated distances to the different classes. According to Ben-Israel *et al.* [45], several relations can be assumed between the distance  $d_k(q^*)$  and its membership probability  $p_k(q^*)$ , including the working principle defined in (15), where  $F(q^*)$  is a function depending only on  $q^*$ .

$$p_k(q^*)e^{d_k(q^*)} = F(q^*). \quad (15)$$

This relation establishes that the probabilities decay exponentially as distances increase. Based on this principle, the membership probability  $p_k(q^*)$  can be defined as in Equation (16), proposed in [45]. Since  $\hat{d}_k(q^*)$  is an estimator of  $d_k(q^*)$ , we use  $d_k(q^*) = \hat{d}_k(q^*)$ . The confidence index RD is the membership probability  $p_k(q^*)$  of the predicted class  $k = y^*$ .

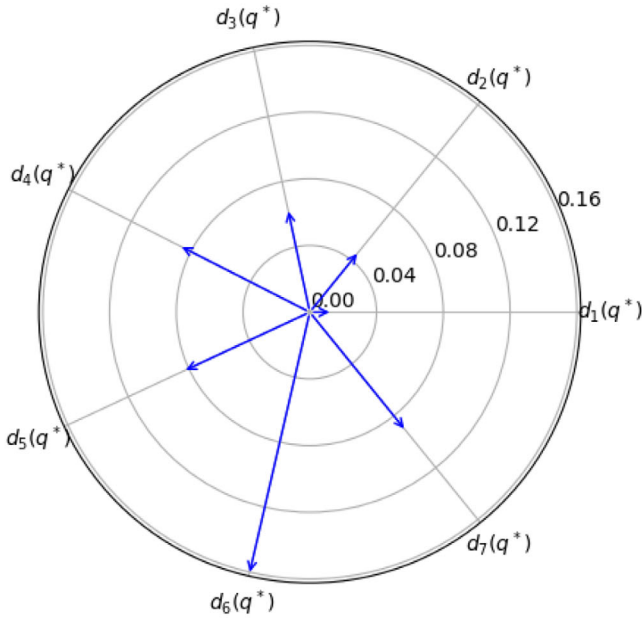
$$p_k(q^*) = \frac{\prod_{j \neq k} e^{d_j(q^*)}}{\sum_{i=1}^K \prod_{j \neq i} e^{d_j(q^*)}}, \quad k = 1, 2, \dots, K. \quad (16)$$

Figure 10 illustrates the distance to class and membership probability of a new voltage sag caused by an upstream symmetrical fault.

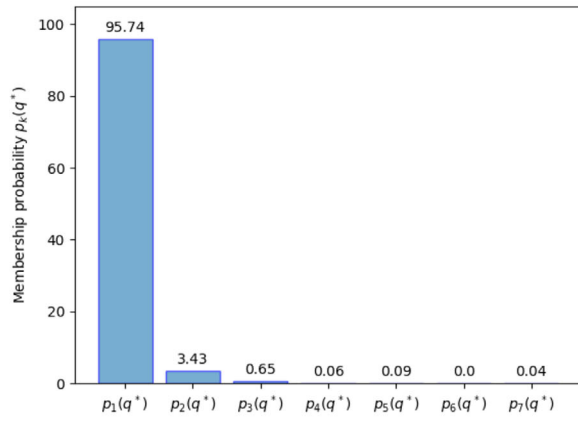
Although an extensive and variate database would naturally improve the estimation of class distances  $d_k(q^*)$ , the relative distance-based index should be less sensitive to the size and diversity of the database than the NB-KDE index.

## 5 | PERFORMANCE ANALYSIS OF THE ALGORITHM

This section investigates the classification performance of the proposed method. For this purpose, a synthetic database is created with a large number of fault and disturbance conditions in order to evaluate the limitations of our proposal and validate the relevance of our methodology.



(a) Distance to class

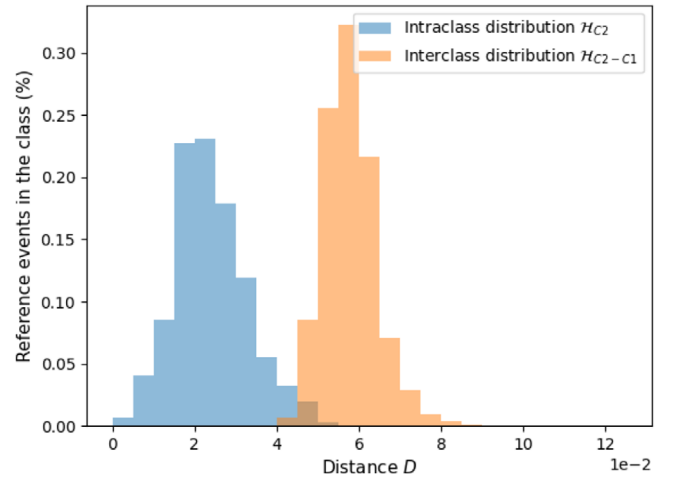


(b) Membership probability

**FIGURE 10** (a) Distances  $d_k(q^*)$  of a new sag  $q^*$  to each class  $k$ , and its corresponding (b) membership probability  $p_k(q^*)$

## 5.1 | Synthetic dataset

A synthetic dataset generated with the simulation model described in section 3 is built. The dataset consists of 700 voltage sags equally distributed in 7 classes (100 events per class). The defined classes are upstream balanced faults (A1), upstream unbalanced faults (A2), downstream balanced faults (B1), downstream unbalanced faults (B2), upstream transformer energizing (C1), downstream transformer energizing (C2), and downstream motor startup (D). The sampling frequency of the synthetic data is set at 12.8 kHz to match the sampling frequency of the real data measurements. However, it can be noted that the sampling frequency could be reduced up to 400 Hz, while maintaining the main characteristics of the four-dimension signatures.



**FIGURE 11** Intra-class distribution  $\mathcal{H}_{C2}$  and inter-class distribution  $\mathcal{H}_{C2-C1}$ , with  $BC = 0.07$

## 5.2 | Class separability

To verify the class separability capabilities of the proposed method, we analyze the characteristics of the entire synthetic dataset. We study the separability of two classes by analyzing their intra-class and inter-class distances. The intra-class distance distribution  $\mathcal{H}_k$  is defined as the ensemble of distances  $D(x_{k,i}, x_{k,j})$  of all the pairs of elements  $x_{k,i}$  and  $x_{k,j}$  belonging to the same class  $C_k$ . The inter-class distance distribution  $\mathcal{H}_{k-k'}$  corresponds to the ensemble of distances  $D(x_{k,i}, x_{k',i})$  between all the pairs of elements  $x_{k,i}$  and  $x_{k',i}$  from two different classes  $C_k$  and  $C_{k'}$ , respectively.

For this study, we use the Bhattacharyya coefficient ( $BC$ ), as defined in (17). This coefficient which varies from 0 to 1, can be interpreted as the overlap between two distributions  $\mathcal{H}_\alpha(x)$  and  $\mathcal{H}_\beta(x)$  defined in  $X$ . Two well-separated classes have a near-to-zero overlap between their intra-class and inter-class distance distributions.

$$BC = \sum_{x \in X} \sqrt{\mathcal{H}_\alpha(x)\mathcal{H}_\beta(x)}. \quad (17)$$

To study the class separability of a particular class  $\alpha$  we calculate the Bhattacharyya coefficient between its intra-class distribution  $\mathcal{H}_\alpha$  and the inter-class distributions to the rest of the classes  $\mathcal{H}_{\alpha-\beta}$ , with  $\beta = \{1, \dots, K\}, \alpha \neq \beta$ . Figure 11 illustrates the intra-class distribution of class C2 and the inter-class distribution between C2 and C1.

For all the considered cases, the Bhattacharyya results are displayed in Table 2. From these results, one can conclude that the cases can be considered sufficiently well separated: all the coefficient values are close to zero (except on the diagonal of the table where the coefficient is computed between 2 identical distributions). The worst Bhattacharyya coefficient is  $BC = 0.07$  obtained for Intra-class distribution  $\mathcal{H}_{C2}$  and inter-class distribution  $\mathcal{H}_{C2-C1}$ .



**TABLE 2** Bhattacharyya coefficient between the intra- and inter-class distributions

Intraclass distribution $\mathcal{H}_\alpha$	Inter-class distribution $\mathcal{H}_{\alpha-\beta}$						
	$\mathcal{H}_{\alpha-A1}$	$\mathcal{H}_{\alpha-A2}$	$\mathcal{H}_{\alpha-B1}$	$\mathcal{H}_{\alpha-B2}$	$\mathcal{H}_{\alpha-C1}$	$\mathcal{H}_{\alpha-C2}$	$\mathcal{H}_{\alpha-D}$
$\mathcal{H}_{A1}$	1.00	0.00	0.00	0.00	0.00	0.00	0.00
$\mathcal{H}_{A2}$	0.01	1.00	0.00	0.00	0.00	0.00	0.00
$\mathcal{H}_{B1}$	0.00	0.00	1.00	0.00	0.00	0.00	0.00
$\mathcal{H}_{B2}$	0.00	0.00	0.04	1.00	0.00	0.00	0.00
$\mathcal{H}_{C1}$	0.03	0.06	0.03	0.04	1.00	0.00	0.00
$\mathcal{H}_{C2}$	0.00	0.00	0.00	0.06	0.07	1.00	0.00
$\mathcal{H}_D$	0.00	0.00	0.04	0.00	0.00	0.00	1.00

### 5.3 | Classification efficiency: Accuracy and robustness

The standard metrics for classification algorithms that are used for evaluating the classification output performances of our algorithm are accuracy, recall, and F1-score (18), the latter being a good metric for summarizing the first two, mainly when applied to balanced datasets. A perfect classification is obtained with an F1-score equal to one.

$$F1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}. \quad (18)$$

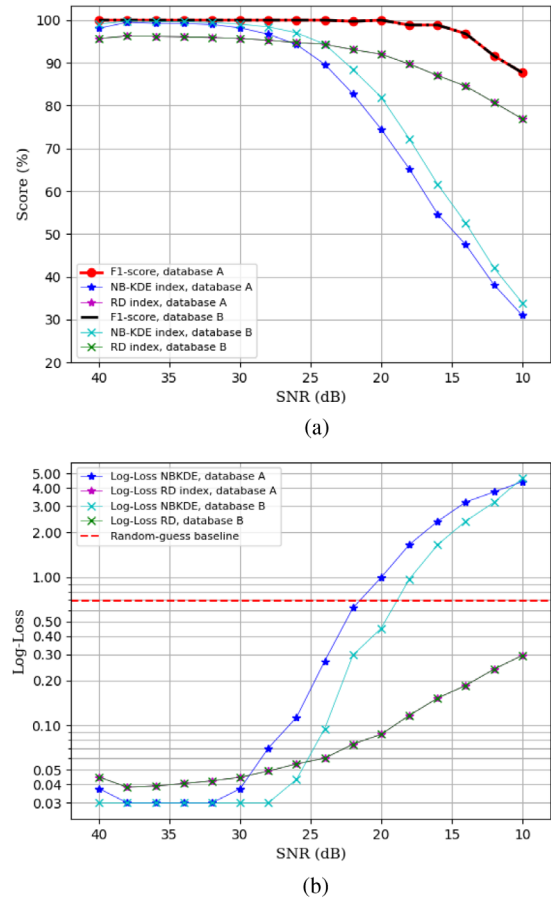
However, the metrics mentioned above are not suitable for evaluating the confidence score associated with the classification output provided by the algorithm. A more appropriate metric for this task is Log-Loss. Globally, this metric penalizes outputs with a low confidence score. The Log-Loss definition is given in (19), where  $y = \{0, 1\}$  and  $p$  is the associated probability estimates with  $p = P(y = 1)$ .

$$\mathcal{L}_{log} = -(y \log(p) + (1 - y) \log(1 - p)). \quad (19)$$

Although raw Log-Loss values can be hard to interpret, lower values mean classification outputs with higher confidence. For instance, a perfect classifier would have a Log-Loss equal to 0. A random-guessing Log-Loss baseline score can be useful for interpreting this metric. The Log-Loss value for  $p = 0.5$  is  $\mathcal{L}_{log} = 0.693$ . Any value higher than this baseline (represented as a red dashed line on the following figures) can be interpreted as worse than random guessing.

We have also studied the robustness to noise of the proposed classifier. For this purpose, we performed two experiments with two databases.

Database A is composed of the original 700 synthetic voltage sags from the synthetic dataset with no added noise. Database B comprises the original 700 synthetic voltage sags plus 700 events with additional white Gaussian noise at SNR=25 dB. To evaluate the performance with increasing noise levels regarding databases A and B, we have created 16 test sets. Each one is composed of the original 700 synthetic voltage sags to which



**FIGURE 12** F1-score and mean values for NB-KDE and RD indexes (a) and the corresponding Log-loss values (b) associated to the classification outputs for different levels of SNR

we add an equal level of additional Gaussian noise, with an SNR varying from 40 to 10 dB.

The results of the experiments with databases A and B are illustrated in Figure 12. F1-score and mean values of the corresponding confidence indexes NB-KDE and RD are displayed in Figure 12a and the Log-loss values for the NB-KDE and RD indexes are shown in Figure 12b.

The results highlight that the evolution of the F1-score is the same for both experiments. It reaches its maximum value between SNR=40 dB and SNR=20 dB, and it slowly decreases in the range from 20 to 10 dB. The confidence indexes NB-KDE and RD maintain a relatively high mean value for all classes from 40 to 20 dB but degrade gradually from 25 to 10 dB. In the same way, their Log-Loss values are low and stable from 40 to 20 dB but increase from 25 to 10 dB.

The NB-KDE index is less robust to variations between the database and the test set, as it rapidly degrades at noise levels higher than SNR=25 dB. We note that the extension of the database with noisy data (database B) slightly improves the performance of the NB-KDE index for SNR levels from 25 to 15 dB. However, beyond 15 dB, the NB-KDE index calculated with database B performs worse with database A, as its Log-Loss values increase. On the contrary, the addition of noisy data

to the database does not seem to influence the RD index, as it is more stable regarding variations in the database.

The global classification performance of our proposal is relatively robust in the standard range of noise levels from 40 to 20 dB. This can be explained by the fact that the pre-processing step implicitly filters the fluctuations through the STFT transform.

## 5.4 | Sensitivity to fundamental frequency variations

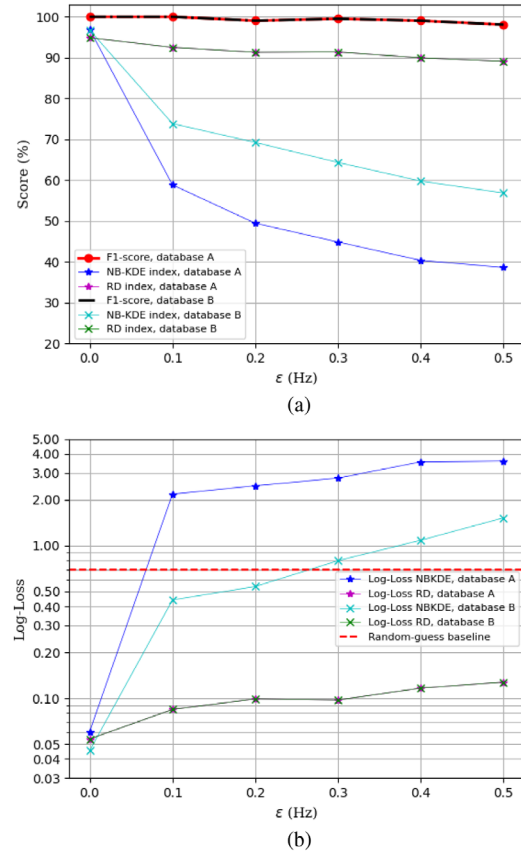
To evaluate the sensitivity of the classification results regarding frequency variations around fundamental frequency (50 Hz), we proceed as before with two experiments based on two databases.

Database A is composed of the original 700 synthetic voltage sags with a frequency of 50.0 Hz, and database B is composed of the 700 synthetic voltage sags plus 210 events at 50.25 Hz (910 events in total). In order to evaluate the performance with increasing levels of frequency variation, we build 6 test sets, each one composed of 210 synthetic voltage sags with a frequency of  $F = 50 \text{ Hz} \pm \epsilon$ , where  $\epsilon = \{0.0, 0.1, 0.2, 0.3, 0.4, 0.5\}$  Hz. This range of frequency variations corresponds to the maximum frequency fluctuations (50 Hz  $\pm 1\%$ ) allowed by French and European regulation standards regarding the power supply at the distribution level for synchronous connection to an interconnected system [46].

The results are plotted in Figure 13a. F1-score values are above 95% for the analyzed range of values for databases A and B. We note that the RD index slowly decreases as  $\epsilon$  increases but stays high near 90%, even for frequency variations of  $\epsilon = \pm 0.5 \text{ Hz}$ . This is true for both databases. We note a significant improvement for the NB-KDE index when the extended database B is used, compared to database A. However, the RD-index is significantly more robust and stable than the NB-KDE index. Regarding Figure 13b, the Log-Loss values of the RD index remain stable despite frequency variations, and are significantly lower than the random-guess baseline for both datasets A and B. This is not the case for the Log-Loss values of the NBKDE index, which increase and perform worse than the random-guess baseline with only a  $\pm 0.1 \text{ Hz}$  frequency variation for database A. Although the performance of the Log-loss values for the NBKDE index slightly improved with database B, its performance is still highly degraded for frequency variations over  $\pm 0.25 \text{ Hz}$ .

## 5.5 | Comparison with benchmark method 1NN-DTW

In this section, we compare the performance of the proposed method with the benchmark method 1NN-DTW presented in [28]. One of the main limitations of 1NN-DTW is its sensitivity to outliers. Indeed, since a single closest neighbor is used for the classification, the presence of outliers or errors in the labeling of the training database can result in classification errors. To evaluate the sensitivity to outliers of both methods, we intro-



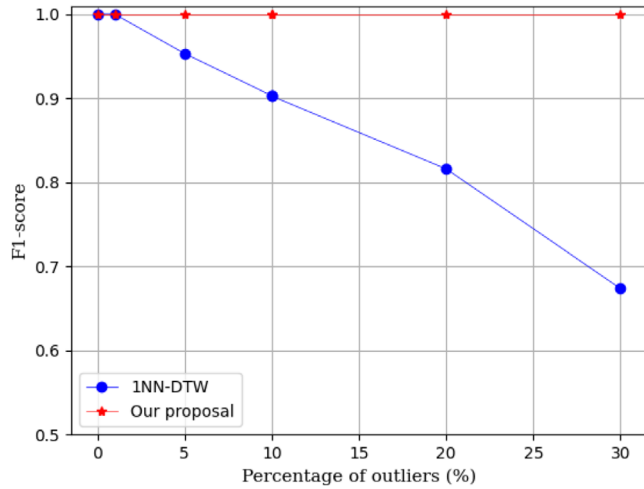
**FIGURE 13** F1-score and mean values for NB-KDE and RD indexes (a) and the corresponding Log-loss values (b) associated to the classification outputs, for different levels of fundamental frequency variations,  $F = 50 \text{ Hz} \pm \epsilon$

duce labeling errors in the reference database using the synthetic dataset, ranging from 0 to 30%. The reference database comprises 60 events per class (420 in total) and the test set comprises 40 events per class (280 in total). The classification results in terms of F1-score are presented in Figure 14 for both methods.

We observe that the new method is more robust to outliers than the 1NN-DTW method. The F1-score for the benchmark method decreases rapidly to 0.95 for 5% of outliers and 0.67 for 30% of outliers. On the contrary, the proposal maintains an F1-score of 1 even with 30% of outliers in the reference database, as the classification is performed taking into account the distance to an entire class instead of a single neighbor. Another significant advantage of the method presented in this work is the provision of confidence scores associated with the classification result, in contrast to 1NN-DTW, which does not provide this information.

## 6 | VALIDATION USING INDUSTRIAL DATA

In order to validate the efficiency of our proposal on real data, the performance on measurement data is presented in this section. The signature database is composed exclusively of



**FIGURE 14** Sensitivity to outliers in the reference (training) database of the proposal compared to 1NN-DTW (benchmark)

**TABLE 3** Real dataset description

Class	Site 1	Site 2	Site 3	Total
A1	-	35	-	35
A2	23	103	29	157
B1	-	-	-	-
B2	-	-	90	90
C1	7	-	2	9
C2	-	-	-	-
D	93	-	1	94
Total	123	138	120	385

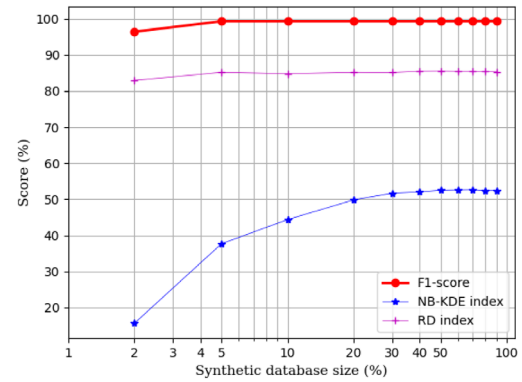
synthetic data collected from the simulation model. The test set contains only real voltage sag records obtained from three industrial sites.

The data consists of records from monitoring devices installed at three sites, each within a specific sector: metal equipment manufacturing, food processing, and chemical industries. The industrial sites under consideration are supplied from the 21 kV distribution network via a 21 kV/400 V *Dyn11* transformer. The voltage drop detection algorithm is integrated into the monitoring device. The dataset consists of 385 measurements, representing 5 of the seven classes defined above, as described in Table 3.

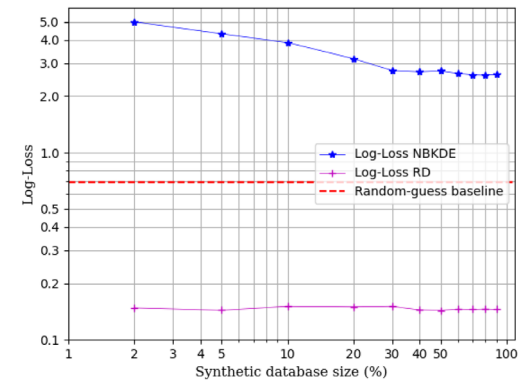
The analysis is performed under the scope of the confusion matrix, accuracy, and recall metrics. The classification errors are also analyzed, and the results highlight the usefulness of the proposed confidence indexes in identifying and avoiding possible false results.

## 6.1 | Minimum database size setting

In the literature, classification studies generally need a huge number of data to perform. Unfortunately in industrial prob-



(a) F1-score and mean values for NB-KDE and RD confidence indexes



(b) Log-loss values for NB-KDE and RD confidence indexes

**FIGURE 15** F1-score and mean values for NB-KDE and RD indexes (a) and the corresponding Log-loss values (b) associated to the classification outputs, for different synthetic database sizes and tested on real data

lems, the availability of data in different faulty conditions is not possible. For this reason, performing the classification with a reduced size is a key point in this area of research. In this section, we determine the minimum database size for synthetic data while maintaining good classification performances. We use a random permutation strategy of five balanced splits. The database is composed of a determined percentage of the synthetic dataset, and the test set is composed at each iteration of the entire experimental dataset. We note that the database contains data from 7 classes, and the test set contains events from the 5 available classes.

The results are illustrated in Figure 15. The F1-score reaches values higher than 95% even with very few samples per class in the database. However, the performance between the confidence indexes is significantly different. The RD-index is noticeably higher and more stable with the increasing number of samples. On the contrary, the NB-KDE index is close to 50% even with a database of maximal size. This poor performance is also reflected in the Log-Loss curve in Figure 15b when compared to the red dashed baseline.

These curves confirm the results obtained in the previous sections. The optimal size of the database is close to 20 samples per class, and the RD-index is the most reliable confidence index compared to the NB-KDE index.

**TABLE 4** Real data results with a synthetic signature database

Class	Database		Test set				
	size	size	Precision	Recall	F1-score	NB-KDE	RD
A1	20	35	94.87	100	97.37	25.58	86.86
A2	20	157	100	98.72	99.36	66.80	92.84
B1	20	-	-	-	-	-	-
B2	20	90	100	100	100	87.96	90.01
C1	20	9	100	100	100	75.24	80.17
C2	20	-	-	-	-	-	-
D	20	94	100	100	100	8.32	74.02
<b>Total</b>	<b>140</b>	<b>385</b>	<b>98.97</b>	<b>99.74</b>	<b>99.34</b>	<b>52.78</b>	<b>84.78</b>

## 6.2 | Results

The reference signature database is composed of 20 synthetic voltage sag events per class or 140 events in total for 7 classes. The synthetic dataset is randomly split into five balanced sets of 140 events; each set is used once as a database. The final results correspond to the average. The test set consists exclusively of experimental measurements.

The obtained results are presented in Table 4. For three out of the five classes, the accuracy, recall, and F1-score are maximal. Only two events of class A2 are misclassified in class A1. The overall results of the classification are satisfactory. However, the mean values for the NB-KDE and RD confidence indexes are lower than those obtained with the synthetic data. This is especially true for the NB-KDE index, which is much lower for classes A1 and D. Nevertheless, the RD index is more stable even when used with field data.

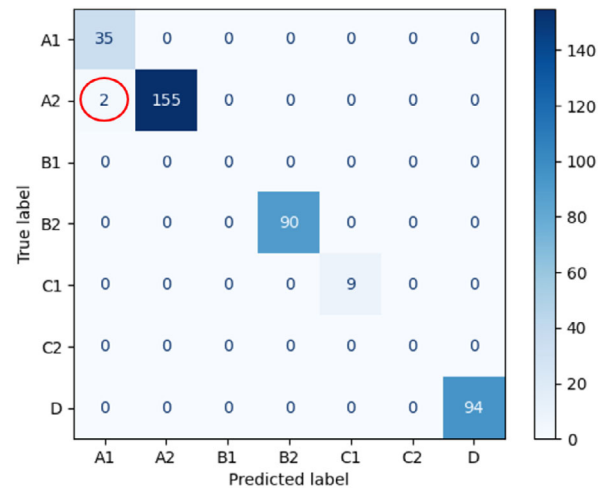
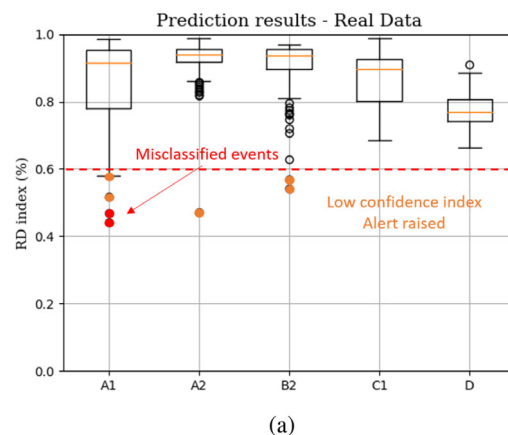
## 6.3 | Accuracy analysis

In this section, we analyze the classification errors leading to the total accuracy of the classifier. Two voltage sags of the real dataset belonging to class A2 were misclassified in class A1, as illustrated in the confusion matrix in Figure 16. If we analyze the relative distance-based confidence index of all the voltage sags, we observe that only seven events out of 385 were classified with a confidence index lower than 60%, as shown in Figure 17. Among these events, two correspond to misclassification errors.

The two events, classified in class A1 have an RD index of 46.89% and 44.11%, respectively. These values are significantly lower than the mean RD index for this class, close to 86%.

In addition to this, the classification error appears between two classes that only differ on the balanced/unbalanced nature of the event. Classes A1 and A2 correspond to voltage sags due to upstream line faults. The consequences of this error could even be seen as minor.

Finally, we can define a threshold at 60% for the RD-index. A classification output with a confidence index below this value would trigger an alert for further analysis by a human expert.

**FIGURE 16** Confusion matrix of results on real data with a synthetic signature database

(a)

True class	Classifier output		RD confidence index (%)						
	Predicted class	RD index	A1	A2	B1	B2	C1	C2	D
A2	A1	46.89	<b>46.89</b>	<b>32.67</b>	1.35	0.81	13.88	0.00	4.40
A2	A1	44.11	<b>44.11</b>	<b>40.35</b>	2.04	1.35	7.90	0.00	4.25
B2	B2	56.84	2.73	<b>33.37</b>	5.44	<b>56.84</b>	0.95	0.00	0.67
B2	B2	54.21	1.61	1.14	<b>34.97</b>	<b>54.21</b>	0.02	0.00	8.05
A1	A1	51.68	51.68	1.02	<b>40.08</b>	0.95	0.47	0.00	5.80
A2	A2	47.14	<b>23.74</b>	<b>47.14</b>	4.81	18.12	1.33	0.01	4.85
A1	A1	58.02	58.02	9.17	10.60	0.67	1.40	0.02	<b>20.12</b>

(b)

**FIGURE 17** Analysis of misclassified events with the real dataset: (a) Plotbox of RD confidence index, 7 events raise an alert due to low confidence index, with two of them corresponding to the misclassification errors. (b) Detailed RD values for the events raising an alert

## 7 | CONCLUSION

This paper has presented a method based on voltage and current waveforms analysis to classify the causes of voltage sags in LV industrial power networks. The proposal is based on four-dimension time series signatures, extracted after a pre-processing stage using STFT and Fortescue transform. A database composed of these signatures and a distance-based classification approach are used for the identification stage.



The performance of the method has been analyzed in terms of class separability, classification efficiency (accuracy and robustness to noise), and sensitivity to fundamental frequency variations. The results have proved that the proposal is resilient regarding noise levels up to 15dB and fundamental frequency variations up to  $\epsilon = \pm 0.5$  Hz. Two confidence indexes, the NB-KDE and RD index, have been proposed and compared. The RD index proved to be more robust and stable. The information provided by such an index increases the reliability on the classification process by alerting when classification results with low confidence scores are obtained while maintaining a high degree of automation in the analysis.

Using a reference database composed entirely of synthetic data, we evaluated the algorithm with synthetic data at first, and finally with field data collected from three industrial sites representing five of the seven defined classes. In both cases, the classification metrics reach high values: 100% for synthetic data and higher than 99% for experimental ones.

The main advantages of the proposed method are (1) the reduced amount of data necessary for building the database and (2) its generalization capabilities. Such an algorithm could be easily implemented in industrial applications with no previous recordings needed since the database can be entirely composed of synthetic data and still provide accurate classification results in different industrial sites regardless of their variety. Other advantages include the (3) electrical interpretability of the signatures and the (4) confidence index associated with the classification output. These characteristics could ease the troubleshooting process and the general interpretability of the algorithm. This is particularly interesting for industrial applications since understanding the algorithm decision-making process is essential for reliability issues.

Future work will include the analysis of voltage sags' effects and consequences on industrial sites. The objective is to estimate the proportion of self-disconnected load types following a voltage sag using the electrical data of a single monitoring device installed at the main MV/LV transformer.

## AUTHOR CONTRIBUTIONS

Maria Veizaga: Conceptualization, data curation, formal analysis, investigation, methodology, resources, software, validation, writing - original draft. Claude Delpha: Conceptualization, formal analysis, funding acquisition, investigation, methodology, supervision, validation, writing - original draft, writing - review and editing. Demba Diallo: Conceptualization, formal analysis, funding acquisition, investigation, methodology, supervision, validation, writing - original draft, writing - review and editing. Sophie Bercu: Conceptualization, investigation, methodology, project administration, supervision, validation, visualization. Ludovic Bertin: Conceptualization, investigation, project administration, supervision, validation, visualization.

## CONFLICT OF INTEREST

The authors have declared no conflict of interest.

## DATA AVAILABILITY STATEMENT

No data

## ORCID

Claude Delpha  <https://orcid.org/0000-0003-3224-8628>

## REFERENCES

1. Targosz, R., Manson, J.: Pan-European power quality survey. In: 2007 9th International Conference on Electrical Power Quality and Utilisation, pp. 1–6. IEEE, Piscataway (2007). Available: <http://ieeexplore.ieee.org/document/4424203/>
2. Wagner, V., Andreshak, A., Staniak, y.: Power quality and factory automation. *IEEE Trans. Ind. Appl.* 26(4), 620–626 (1990). Available: <http://ieeexplore.ieee.org/document/55984/>
3. Sarmiento, H., Estrada, E.: A voltage sag study in an industry with adjustable speed drives. *IEEE Ind. Appl. Mag.* 2(1), 16–19 (1996)
4. IEEE: IEEE Recommended Practice for Monitoring Electric Power Quality. *IEEE Std 1159-2019* (Revision of IEEE Std 1159-2009), pp. 1–98 (2019)
5. IEC: IEC TR 61000-2-8:2002 - Voltage dips and short interruptions on public electric power supply systems with statistical measurement results, Tech. Rep. IEC/CEI 61000-2-8:2002 (2002)
6. Bollen, M.H.J.: *Understanding Power Quality Problems*. IEEE Press, New York (2000)
7. Bollen, M., Zhang, L.: Different methods for classification of three-phase unbalanced voltage dips due to faults. *Electr. Power Syst. Res.* 66(1), 59–69 (2003). Available: <https://linkinghub.elsevier.com/retrieve/pii/S0378779603000725>
8. Ignatova, V., Granjon, P., Bacha, S.: Space vector method for voltage dips and swells analysis. *IEEE Trans. Power Delivery* 24(4), 2054–2061 (2009). Available: <http://ieeexplore.ieee.org/document/5235772/>
9. Thakur, P., Singh, A.K.: Unbalance voltage sag fault-type characterization algorithm for recorded waveform. *IEEE Trans. Power Delivery* 28(2), 1007–1014 (2013). Available: <http://ieeexplore.ieee.org/document/6407157/>
10. Bollen, M.H.J., Gu, I.Y., Axelberg, P.G., Styvaktakis, E.: Classification of underlying causes of power quality disturbances: Deterministic versus statistical methods. *EURASIP J. Adv. Signal Process.* 2007(1), 079747 (2007). Available: <https://asp-urasipjournals.springeropen.com/articles/10.1155/2007/79747>
11. Erişti, H., Yıldırım, O., Erişti, B., Demir, Y.: Automatic recognition system of underlying causes of power quality disturbances based on s-transform and extreme learning machine. *Int. J. Electr. Power Energy Syst.* 61, 553–562 (2014). Available: <https://linkinghub.elsevier.com/retrieve/pii/S0142061514001938>
12. Mishra, M., Panigrahi, R.R., Advanced signal processing and machine learning techniques for voltage sag causes detection in an electric power system. *Int. Trans. Electr. Energy Syst.* 30(1), (2020). Available: <https://onlinelibrary.wiley.com/doi/10.1002/2050-7038.12167>
13. Manjula, M., Mishra, S., Sarma, A.: Empirical mode decomposition with hilbert transform for classification of voltage sag causes using probabilistic neural network. *Int. J. Electr. Power Energy Syst.* 44(1), 597–603 (2013). Available: <https://linkinghub.elsevier.com/retrieve/pii/S0142061512004000>
14. Wang, N., Wang, S., Jia, Q.: The method to reduce identification feature of different voltage sag disturbance source based on principal component analysis. In: 2014 IEEE Conference and Expo Transportation Electrification Asia-Pacific (ITEC Asia-Pacific), pp. 1–6. IEEE, Piscataway (2014). Available: <https://ieeexplore.ieee.org/document/6940964>
15. Liu, J., Song, H., Zhou, L.: Identification and location of voltage sag sources based on multi-label random forest. In: 2019 IEEE Sustainable Power and Energy Conference (ISPEC), pp. 2025–2030. IEEE, Piscataway (2019). Available: <https://ieeexplore.ieee.org/document/8975097/>
16. Aggarwal, A., Saini, M.K.: Recognition of voltage sag causes using vector quantization based orthogonal wavelet. In: 2020 IEEE 9th Power India

- International Conference (PIICON), pp. 1–6. IEEE, Piscataway (2020). Available: <https://ieeexplore.ieee.org/document/9112953/>
17. Saini, M.K., Aggarwal, A.: Fractionally delayed Legendre wavelet transform based detection and optimal features based classification of voltage sag causes. *J. Renew. Sustain. Energy* 11(1), 015503 (2019). Available: <http://aip.scitation.org/doi/10.1063/1.5049189>
  18. Nagata, E.A., Ferreira, D.D., Bollen, M.H., Barbosa, B.H., Ribeiro, E.G., Duque, C.A., Ribeiro, P.F.: Real-time voltage sag detection and classification for power quality diagnostics. *Measurement* 164, 108097 (2020). Available: <https://linkinghub.elsevier.com/retrieve/pii/S0263224120306357>
  19. Ben-David, S., Blitzer, J., Crammer, K., Pereira, F.: Analysis of Representations for Domain Adaptation. In: Schölkopf, B., Platt, J., Hoffman, T. (eds.) *Advances in Neural Information Processing Systems*, Vol. 19. MIT Press, Cambridge (2006). Available: <https://proceedings.neurips.cc/paper/2006/file/b1b0432ceafb0ce714426e9114852ac7-Paper.pdf>
  20. Susto, G.A., Cenedese, A., Terzi, M.: Chapter 9 - Time-series classification methods: Review and applications to power systems data. In: R. Arghandeh, Y. Zhou (eds.) *Big Data Application in Power Systems*, pp. 179–220. Elsevier, Amsterdam, Netherlands (2018). Available: <https://www.sciencedirect.com/science/article/pii/B9780128119686000097>
  21. Bagnall, A., Bostrom, A., Large, J., Lines, J.: The great time series classification bake off: an experimental evaluation of recently proposed algorithms. extended version. *Data Min. Knowl. Discovery* 31, 606–660 (2017). arXiv: 1602.01711. Available: <http://arxiv.org/abs/1602.01711>
  22. Ruiz, A.P., Flynn, M., Large, J., Middlehurst, M., Bagnall, A.: The great multivariate time series classification bake off: A review and experimental evaluation of recent algorithmic advances. *Data Min. Knowl. Discovery* 35(2), 401–449 (2021). Available: <http://link.springer.com/10.1007/s10618-020-00727-3>
  23. Li, D., Mei, F., Zhang, C., Sha, H., Zheng, J.: Self-Supervised voltage sag source identification method based on CNN. *Energies* 12(6), 1059 (2019). Available: <https://www.mdpi.com/1996-1073/12/6/1059>
  24. Zheng, Z., Qi, L., Wang, H., Zhu, M., Chen, Q.: Recognition method of voltage sag causes based on Bi-LSTM. *IEEE Trans. Electr. Electron. Eng.* 15(3), 418–425 (2020). Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/tee.23070>
  25. Wang, H., Qi, L., Ma, Y., Jia, J., Zheng, Z.: Method of voltage sag causes based on bidirectional lstm and attention mechanism. *J. Electr. Eng. Technol.* 15(3), 1115–1125 (2020). Available: <http://link.springer.com/10.1007/s42835-020-00413-v>
  26. Sakoe, H., Chiba, S.: Dynamic programming algorithm optimization for spoken word recognition. *IEEE Trans. Acoust. Speech Signal Process.* 26(1), 43–49 (1978). Available: <http://ieeexplore.ieee.org/document/1163055/>
  27. Youssef, A., Abdel-Galil, T., El-Saadany, E., Salama, M.: Disturbance classification utilizing dynamic time warping classifier. *IEEE Trans. Power Delivery* 19(1), 272–278 (2004). Available: <http://ieeexplore.ieee.org/document/1256388/>
  28. Veizaga, M., Bercu, S., Delpha, C., Diallo, D., Bertin, L.: Classification of voltage sag causes based on instantaneous symmetrical components using 1NN and dynamic time warping. In: *IECON 2021–47th Annual Conference of the IEEE Industrial Electronics Society*, pp. 1–6. IEEE, Piscataway (2021). Available: <https://ieeexplore.ieee.org/document/9589719/>
  29. Barrera Nunez, V., Gu, I.Y.-H., Bollen, M.H., Melendez, J.: Feature characterization of power quality events according to their underlying causes. In: *Proceedings of 14th International Conference on Harmonics and Quality of Power - ICHQP 2010*, pp. 1–8. IEEE, Piscataway (2010). Available: <http://ieeexplore.ieee.org/document/5625496/>
  30. Ding, N., Cai, W., Suo, J., Wang, J., Xu, Y.: Voltage sag disturbance detection based on RMS voltage method. In: *2009 Asia-Pacific Power and Energy Engineering Conference*, pp. 1–4. IEEE, Piscataway (2009). Available: <http://ieeexplore.ieee.org/document/4918960/>
  31. Alam, M.R., Muttaqi, K.M., Bouzardoum, A.: Characterizing voltage sags and swells using three-phase voltage ellipse parameters. *IEEE Trans. Ind. Appl.* 51(4), 2780–2790 (2015). Available: <http://ieeexplore.ieee.org/document/7027818/>
  32. Camarillo-Penaranda, J.R., Ramos, G.: Fault classification and voltage sag parameter computation using voltage ellipses. *IEEE Trans. Ind. Appl.* 55(1), 92–97 (2019). Available: <https://ieeexplore.ieee.org/document/8428464/>
  33. Yalçinkaya, G., Bollen, M., Crossley, P.: Characterization of voltage sags in industrial distribution systems. *IEEE Trans. Ind. Appl.* 34(4), 682–688 (1998). Available: <http://ieeexplore.ieee.org/document/703958/>
  34. Oubrahim, Z., Choqueuse, V., Amirat, Y., Benbouzid, M.E.H.: Disturbances classification based on a model order selection method for power quality monitoring. *IEEE Trans. Ind. Electron.* 64(12), 9421–9432 (2017). Available: <http://ieeexplore.ieee.org/document/7938438/>
  35. Polajžer, B., Štumberger, G., Dolinar, D.: Instantaneous positive-sequence current applied for detecting voltage sag sources. *IET Gener. Transm. Distrib.* 9(4), 319–327 (2015). Available: <https://onlinelibrary.wiley.com/doi/10.1049/iet-gtd.2014.0483>
  36. Mohammadi, Y., Moradi, M.H., Chouhy Leborgne, R.: Employing instantaneous positive sequence symmetrical components for voltage sag source relative location. *Electr. Power Syst. Res.* 151, 186–196 (2017). Available: <https://linkinghub.elsevier.com/retrieve/pii/S0378779617302286>
  37. Bollen, M.H.J., Gu, I.Y.-H.: *Signal Processing of Power Quality Disturbances*. IEEE Press Series on Power Engineering, Wiley-Interscience, Hoboken, NJ (2006)
  38. Mahseredjian, J., Denetiere, S., Dubé, L., Khodabakhchian, B., Gérin-Lajoie, L.: On a new approach for the simulation of transients in power systems. *Electr. Power Syst. Res.* 77(11), 1514–1520 (2007). Available: <https://www.sciencedirect.com/science/article/pii/S0378779606002094>
  39. Kehtarnavaz, N.: Chapter 7 - Frequency domain processing. In: Kehtarnavaz, N. (ed) *Digital Signal Processing System Design*, pp. 175–196, 2nd ed. Academic Press, Burlington (2008). Available: <https://www.sciencedirect.com/science/article/pii/B9780123744906000076>
  40. Fortescue, C.L.: Method of symmetrical co-ordinates applied to the solution of polyphase networks. *Trans. Am. Inst. Electr. Eng.* XXXVII(2), 1027–1140 (1918). Available: <http://ieeexplore.ieee.org/document/4765570/>
  41. Aung, M., Milanovic, J.: The influence of transformer winding connections on the propagation of voltage sags. *IEEE Trans. Power Delivery* 21(1), 262–269 (2006). Available: <http://ieeexplore.ieee.org/document/1564208/>
  42. Parsons, A., Grady, W., Powers, E., Soward, J.: A direction finder for power quality disturbances based upon disturbance power and energy. *IEEE Trans. Power Delivery* 15(3), 1081–1086 (2000). Available: <http://ieeexplore.ieee.org/document/871378/>
  43. Itakura, F.: Minimum prediction residual principle applied to speech recognition. *IEEE Trans. Acoust. Speech Signal Process.* 23(1), 67–72 (1975). Available: <http://ieeexplore.ieee.org/document/1162641/>
  44. Giorgino, T.: Computing and visualizing dynamic time warping alignments in R: The dtw package. *J. Stat. Softw.* 31(7), 1–24 (2009). Available: <http://www.jstatsoft.org/v31/i07/>
  45. Ben-Israel, A., Iyigun, C.: Probabilistic D-clustering. *J. Classif.* 25(1), 5 (2008). Available: <https://doi.org/10.1007/s00357-008-9002-z>
  46. CEN-CENELEC: NF EN 50160:2011 - Voltage characteristics of electricity supplied by public electricity networks (2011)

**How to cite this article:** Veizaga, M., Delpha, C., Diallo, D., Bercu, S., Bertin, L.: Classification of voltage sags causes in industrial power networks using multivariate time-series. *IET Gener. Transm. Distrib.* 1–17 (2023). <https://doi.org/10.1049/gtd.2.12765>