

Multi-Class Disturbance Events Recognition Based on EMD and XGBoost in φ -OTDR

ZHANDONG WANG¹, SHUQIN LOU¹, SHENG LIANG², AND XINZHI SHENG²

¹School of Electronic and Information Engineering, Beijing Jiaotong University, Beijing 100044, China

²School of Science, Beijing Jiaotong University, Beijing 100044, China

Corresponding author: Xinzhi Sheng (xzhsheng@bjtu.edu.cn)

This work was supported in part by the National Natural Science Foundation of China under Grant 61775014.

ABSTRACT A novel pattern recognition method based on Empirical Mode Decomposition (EMD) and extreme gradient boosting (XGBoost) is proposed to recognize the disturbance events in phase sensitive optical time-domain reflectometer (φ -OTDR) to reduce nuisance alarm rate (NAR) and improve real-time performance in this paper. Eleven typical eigenvectors are extracted from components obtained by EMD of the disturbance signals and XGBoost is selected as a classifier to identify different type of disturbance signals. Five kinds of disturbance events, including watering, knocking, climbing, pressing and false disturbance event, can be identified, effectively. Experimental results show that NAR is 4.10% and identification time is 0.093 s. The recognition accuracy for the five patterns is 97.96%, 95.90%, 91.10%, 94.84% and 99.69%, respectively. The effectiveness of the proposed method is evaluated by using confusion matrix and decision boundary visualization. Experimental results demonstrate that our proposed pattern recognition method based on XGBoost has better performance in recognition rate and recognition time than other commonly used methods, such as support vector machine (SVM), Gradient Boosting Decision Tree (GBDT), Random Forest (RF) and Adaptive Boosting (Adaboost).

INDEX TERMS Phase-sensitive optical time-domain reflectometer (φ -OTDR), extreme gradient boosting (XGBoost), nuisance alarm rate (NAR), empirical mode decomposition (EMD), pattern recognition, decision boundary visualization.

I. INTRODUCTION

Recently, fiber-optic distributed disturbance sensors based on phase-sensitive optical time-domain reflectometer (φ -OTDR) have drawn intensive attention due to their low loss, simple structure, chemical stability, anti-electromagnetic interference. Due to advantages of high spatial resolution, high sensitivity, long-distance transmission capability and accurate location of disturbance, φ -OTDR has a potential application in the field of perimeter security [1], [2], speed monitoring and localization of trains [3], pipeline security [4], [5], and vibration measurement [6]–[8].

However, in the practical φ -OTDR monitoring system, external environment interference and artificial non-destructive interference often cause false alarms, resulting in a high nuisance alarm rate (NAR) [9]. Generally, there are two strategies used to reduce NAR. One is to improve

the hardware structure and performance [10]. For instance, the monitoring system adopt a combined structure of φ -OTDR and fiber-optic interferometers i.e. Mach-Zehnder Interference [11] or Michelson Interferometers [12]. These schemes can efficiently reduce NAR, but hardware structure become more complicated and thus the system cost increases significantly.

Another is to select and optimize pattern recognition methods to improve the recognition accuracy of disturbance events. In 2015, a method based on morphological feature extraction of time-space domain signals and relevance vector machine (RVM) classifier was reported [13]. The average identification rate of three events reached up to 97.8%, but the recognition time was 0.7 s. In 2017, Xu *et al.* [14] used spectral subtraction to reduce wide-band background noise of signals and support vector machine (SVM) to detect four disturbance events (taping, striking, shaking, and crushing). The average identification rate was 93.8% and identification time was below 0.6 s. Subsequently, they reported

The associate editor coordinating the review of this manuscript and approving it for publication was Cesar Vargas-Rosales ¹.

a pattern recognition method based on convolution neural network (CNN) and SVM in 2018 [15]. Through replacing the soft-max classifier in CNN with a nonlinear SVM classifier or a linear SVM classifier, the identification rate exceeded 93.3%. In 2019, Wang et al. [16] used relevant vector machine (RVM) based on a 7-dimensional feature vector extracted by wavelet energy spectrum analysis to identify three disturbance events (walking through the fiber, striking on the fiber, and jogging along the fiber). A classification macro-accuracy of 88.60% was finally obtained through 10-fold cross validation. But the accuracy of walking through the fiber and jogging along the fiber were both under 85%. In the same year, Shi et al. [17] used the data matrix obtained by bandpass filtering and grayscale conversion preprocessing of the original spatiotemporal signal as the input of CNN. By simplifying the structure of GoogLeNet, the running memory was reduced and running speed was improved. The recognition rate of five types of events reached 96.67%. It can be seen that a higher identification rate can be achieved through the selection and optimization of feature extraction and classifier. However, it is still a challenge to find the optimal features and high-performance classifier to further improve the identification rate and reduce the identification time.

Considering non-stationarity features of the disturbance signal in φ -OTDR, EMD is often used to deal with non-stationarity signals [18]. By EMD analysis, the disturbance signals can be decomposed layer by layer according to the frequency characteristics and then the feature vector can be extracted. Meanwhile, the XGBoost algorithm can be used to reduce the identification time as it converges fast and supports parallel operations.

In this paper, we propose a novel pattern recognition method based on the combination of EMD energy analysis and XGBoost algorithm to reduce NAR and shorten recognition time of disturbance events in φ -OTDR. We firstly explain the pattern recognition system based on φ -OTDR briefly, and then describe the basic principles of EMD energy analysis and XGBoost used in this system. The feature vectors are determined by EMD energy analysis. Finally, the XGBoost will be used to identify the disturbance signals. In order to verify the effectiveness of the proposed method, we set up an φ -OTDR experimental system with the length of 25.05 km sensing fiber. Five kinds of disturbance events, including watering, knocking, climbing, pressing, and false disturbance event (false), are considered for identification. We evaluate the classification model by using confusion matrix and performance parameters, and use decision boundaries to compare the classification performance with other methods by visualizing the classification in a two-dimensional feature space. In the rest of the paper, we introduce the pattern recognition method based on EMD energy analysis and XGBoost in Section II. The design of experiment system is described in Section III. The experimental results and discussion are presented in Section IV, and Section V provides some conclusions and ideas for future work.

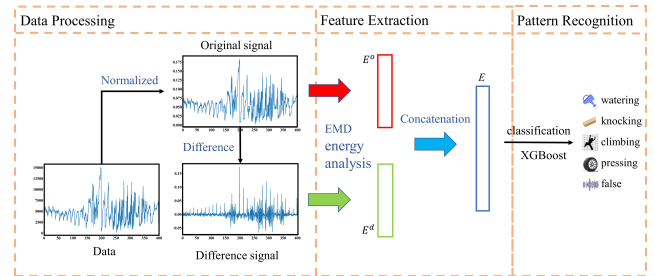


FIGURE 1. The flow chart of a new pattern recognition method.

II. PARRERN RECOGNITION BASED ON XGBoost ALGORITHM

In this paper, we propose a pattern recognition method with the combination of EMD and XGBoost in φ -OTDR. The proposed pattern recognition method includes three steps: data processing, feature extraction and pattern recognition of disturbance events in φ -OTDR, as shown in Fig. 1. The first step is to acquire and process the data from φ -OTDR. The original signal can be obtained by normalizing the data collected by data acquisition card in the experimental system to eliminate the effects of the data itself, and then its difference signal can be obtained by first-order difference processing. The second step is to extract feature vector of the original normalized data and differential signal by using EMD energy analysis. Then XGBoost is chosen as the classifier to recognition five types of signals in pattern recognition step.

A. EMPIRICIAL MODE DECOMPOSITION

EMD is an adaptive decomposition algorithm for processing signals [19]. Compared with the wavelet decomposition method [20], it does not need to set the basis function in advance. The signal can be decomposed into a finite number of the time-domain Intrinsic Mode Functions (IMFs) according to the time scale characteristics of the data itself. The flow chart of EMD energy analysis shown in Fig.2.

According to the EMD theory, the signal $X(n)$ can be expressed as the sum of all IMF components $IMF_l(n)$ and residual component, $res_L(n)$.

$$X(n) = \sum_{l=1}^L IMF_l(n) + res_L(n) \quad (1)$$

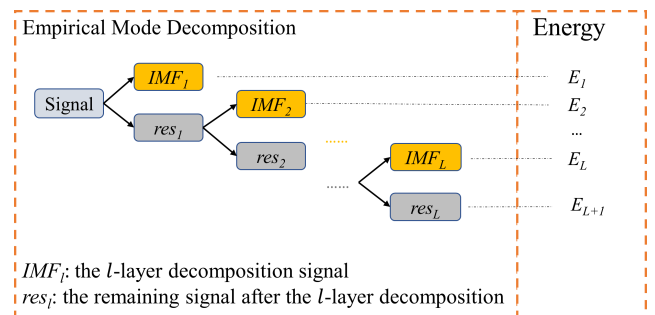


FIGURE 2. The flow chart of EMD energy analysis.

where $IMF_l(n)$ is the l -layer decomposition signal, L is the total number of decomposition levels, and $res_L(n)$ is the remaining signal after L -layer decomposition.

The short-time energy can reflect the change of different signal energies. The short-term energy value E of each component can be calculated as follows.

$$\begin{cases} E_l = \sum_{n=1}^N IMF_l^2(n) \\ E_{L+1} = \sum_{n=1}^N res_L^2(n), \end{cases} \quad l = 1, 2, \dots, L \quad (2)$$

where E_l represents the energy of IMF_l , E_{L+1} represents the energy of res and N represents the length of components of the signal.

Two feature vectors E^o and E^d can be obtained by performing EMD energy analysis on the original signal and differential signal, which can be expressed as.

$$E^o = [E_1^o, E_2^o, \dots, E_{L^o}^o, E_{L^o+1}^o] \quad (3)$$

$$E^d = [E_1^d, E_2^d, \dots, E_{L^d}^d, E_{L^d+1}^d] \quad (4)$$

where L^o and L^d are the decomposition levels of the original signal and differential signal, respectively.

The final feature vector $E = [E^o, E^d]$ is obtained by cascading E^o and E^d , and the number of features is $L^o + L^d + 2$.

B. EXTREME GRADIENT BOOSTING

The choice and design of the classifier is an important part in the pattern recognition since it can improve the classification accuracy and reduce recognition time.

The XGBoost algorithm [21] is improved and introduced as a robust decision tree by Chen based on the idea of a gradient boosting machine [22], which can handle complex data at high speed and accuracy. The feature E in φ -OTDR system and their corresponding signal types constitute a data set $D = \{E_i, y_i\}_{i=1}^N$, where y_i is the type of signal corresponding to feature E_i and N is the number of signal types. K additive functions will be used to predict the output values of a tree ensemble model as follows.

$$\hat{y}_i = \sum_{k=1}^K f_k(E_i), \quad f_k \in F \quad (5)$$

where \hat{y}_i is predicted values and F is the regression trees space which is expressed as:

$$F = \{f(E) = W_q(E)\}(q: R^m \rightarrow T, w \in R^T) \quad (6)$$

where q denotes for the structure of each tree; T denotes for the number of leaves in the tree, and f_k is a function that corresponds to an independent tree structure q and leave weight w .

To reduce errors of ensemble trees, the objective function is calculated in the XGBoost model as below.

$$L^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{t-1} + f_i(E_i)) + \Omega(f_i) \quad (7)$$

where l is a differentiable convex objective function to determine the error between predicted and measured values, t denotes the repetitions in order to minimize the errors and Ω is regression tree complexity penalty function:

$$\Omega(f_k) = \gamma T + \frac{1}{2} \lambda \|w\|^2 \quad (8)$$

where γ and λ are the penalty coefficients.

In this paper, for the best performance of XGBoost, K in equation (5) is 70. The `max_depth`, `colsample_bytree` and `min_child_weight` of structure q in equation (6) are 4, 0.8 and 2.7, respectively [23], [24]. CART tree [25] is chosen as the regression trees F in equation (6). In equation (7), softmax and Mean Squared Error (MSE) are chosen as the objective function L and loss function l , respectively. γ and λ in equation (8) is 0 and 1, respectively.

C. THE CLASSIFIER PERFORMANCE

To evaluate the effectiveness of classifier, the concept of confusion matrix and decision boundary are introduced.

Confusion matrix is a concept from machine learning, which contains information about actual and predicted classifications given by a classification system [26]. A confusion matrix has two-dimensions, in which one dimension is indexed by the actual class of an object and the other is indexed by the class that the classifier predicts. Fig.3 presents the basic form of confusion matrix with the classes A_1, A_2 , and A_n for a multi-classification task. In the confusion matrix, N_{ij} represents the number of samples actually belonging to class A_i but being classified as class A_j .

A number of measures of classification performance can be defined based on the confusion matrix. Some common measures include accuracy, precision, recall and traditional F-score [27].

Accuracy (Acc) is the proportion of the correctly classified sample number to all sample number.

$$Acc = \frac{\sum_{i=1}^n N_{ii}}{\sum_{i=1}^n \sum_{j=1}^n N_{ij}} \quad (9)$$

| | | Predicted | | | | |
|--------|----------|-----------|-----|----------|-----|----------|
| | | A_1 | ... | A_j | ... | A_n |
| Actual | A_1 | N_{11} | | N_{1j} | | N_{1n} |
| | \vdots | | | \vdots | | |
| | A_i | N_{i1} | ... | N_{ij} | ... | N_{in} |
| | \vdots | | | \vdots | | |
| | A_n | N_{n1} | | N_{nj} | | N_{nn} |

FIGURE 3. Confusion matrix representation.

Precision (P) is a measure of the accuracy provided that a specific class has been predicted. It is defined as

$$P_i = N_{ii} / \sum_{k=1}^n N_{ki} \quad (10)$$

Recall (R) is a measure of the ability of a prediction model to select instances of a certain class from a dataset. It can be expressed as

$$R_i = N_{ii} / \sum_{k=1}^n N_{ik} \quad (11)$$

The traditional F-score (F) is the harmonic mean of precision and recall, which can be expressed as

$$F_i = \frac{2 \times P_i \times R_i}{P_i + R_i} \quad (12)$$

A class can also be separated from other classes by visualizing the decision boundary in a classification task [28]. In this paper, we will choose the two most important features according to the importance of the features to build a two-dimensional feature space. By visualizing decision boundaries, we can intuitively see the distribution of different patterns in the feature space and the classification effect of different classifiers such as SVM [29], GBDT [30], RF [31] and Adaboost [32].

III. SYSTEM

A. EXPERIMENT SETUP

The structure of φ -OTDR disturbance sensing system is shown in Fig.4. This system uses a 1550nm continuous laser as light source with a line width of 10 kHz, which is then modulated into a laser pulse sequence with a pulse width of 500 ns by the acousto-optic modulator(AOM). The optic pulse is amplified by an erbium-doped fiber amplifier (EDFA) and injected through a circulator into a single-mode sensing fiber with a total length of 25.05 km. The Rayleigh backscattered light goes back to a photodiode(PD) via the circulator, which is then pre-amplified by an electronic amplifier (AMP) and converted to a digital signal by an analog-to-digital converter (ADC). Finally, the digital signal is collected into computer to be processed and classified. Meanwhile, the disturbance alarm and location can be obtained in real time.

There are five kinds of real disturbance events involving watering (pouring water perpendicularly from a container with a capacity of 300 mL, 50 cm above sensing fiber to simulate rain environment), knocking (knocking the sensing

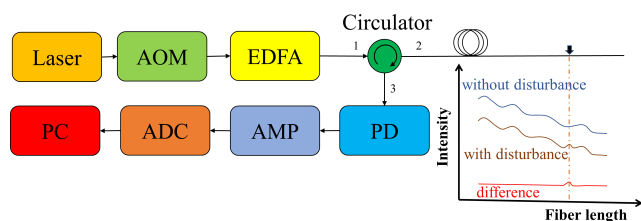


FIGURE 4. The structure of φ -OTDR disturbance sensing system.

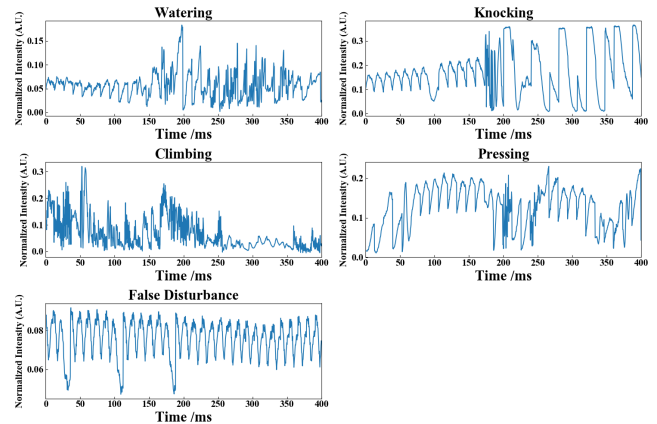


FIGURE 5. The temporal waveforms of normalized signals for five kind of events at the location of 0.9 km.

fiber by using a rod with a radius of 2 cm to simulate stress damage), pressing (using a tire with a radius of 30 cm rolled over the sensing fiber to simulate non-destructive human disturbance such as vehicle passing), climbing (climbing the fence with sensing fiber to simulate destructive human disturbance) and a false disturbance event (induced by the noises which may lead to false alarm) in our experiment. We set up three disturbance locations of 0.9 km, 6.5 km and 21.5 km in the experiment and collect the signals of five disturbance events in each disturbance location simultaneously.

B. SIGNAL PROCESSING AND GROUPING

Data normalization is crucial for feature extraction and event identification for its ability to eliminate the effects of the data itself and overcome the over-fitting problem. Hence, the original signals are normalized in space-domain by L_2 normalization. The normalized time-domain waveforms of the five event patterns at the location of 0.9 km are shown in Fig.5. Moreover, we perform a first-order differential processing on the signal because the differential signal can reflect the changing characteristics of the original signal for further signal analysis and feature extraction. The first-order differential signals are shown in Fig.6.

Before extracting feature values, the signal needs to be divided into several data samples with length N . N is set as 33 sampling periods (13.2ms) according to the temporal period of the error disturbance signal. To ensure the continuity of the signal and increase the number of samples, it requires partial overlap between adjacent data samples. Since each data sample must have its own separate portion, the overlap length is set to one-third of the sample length (4.4ms). The numbers of samples for five kinds of events at each location is shown in Table 1.

IV. EXPERIMENTAL RESULTS AND DISCUSSIONS

A. DETERMINATION OF IMFS FOR FEATURE EXTRACTION

In order to improve the accuracy of recognition and reduce NAR, it is necessary to determine the number of

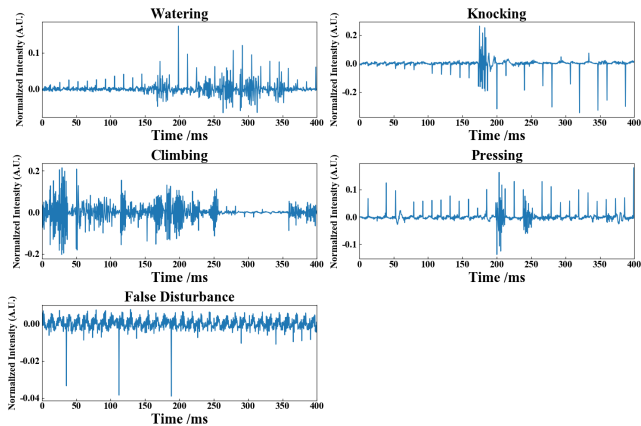


FIGURE 6. Intensity of the first-order differential for five kind of events at the location of 0.9 km.

TABLE 1. The numbers of samples for five kinds of events at each location.

| Event | Disturbance location (km) | | | Total |
|-------|---------------------------|----------------------------------|------|-------|
| | 0.9 | 6.5 | 21.5 | |
| w | 50 | 50 | 50 | 150 |
| k | 50 | 50 | 50 | 150 |
| c | 50 | 50 | 50 | 150 |
| p | 50 | 50 | 50 | 150 |
| f | 50 | 50 </td <td>50</td> <td>150</td> | 50 | 150 |
| Total | 250 | 250 | 250 | 750 |

(Note: w: watering, k: knocking, c: climbing, p: pressing, f: false)

decomposition layers of EMD. In this part, 250 samples from five kinds of disturbance events at the location of 0.9 km are obtained, 60% (150 samples) of the whole samples are divided randomly for training samples and the other 40% (100 samples) for testing samples.

Fig.7 shows the average accuracy of XGBoost based on different decomposition levels for original and differential signals. It can be found that the recognition accuracy increases with the increase of decomposition layers for both original signal and differential signal. The accuracy is approximate to a constant until the decomposition levels exceed 5 layers for the original signal and 4 layers for differential signal. In fact, the more the decomposition levels are, the higher the recognition accuracy is. But the time for feature extraction also increases with the increase of decomposition layers. Considering the tradeoff of recognition time and recognition accuracy, we choose five decomposition levels for the original signal and four for differential signal. So, L^o is 5 and L^d is 4 in the following test, and final number of features is set to 11.

According to separation degree of the 11 features, we select one of the six original signal features E_o^6 and one of the five differential signal characteristics E_d^1 in which the five disturbance events have obvious distinctions, and the selected 2 features are shown in Fig.8. As shown, the point and its color in Fig.8 represent the signal sample and its corresponding category, respectively.

The feature extraction time and average recognition accuracy at the location of 0.9 km of different methods are shown in Table 2.

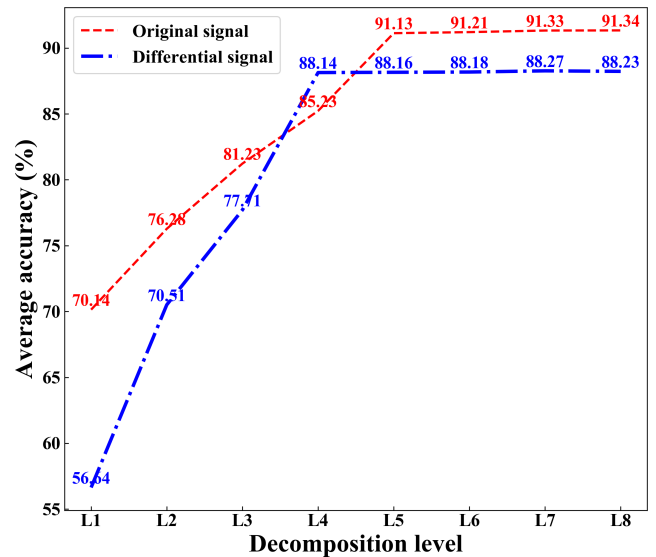


FIGURE 7. The average accuracy of XGBoost based on different decomposition levels for the original and differential signals.

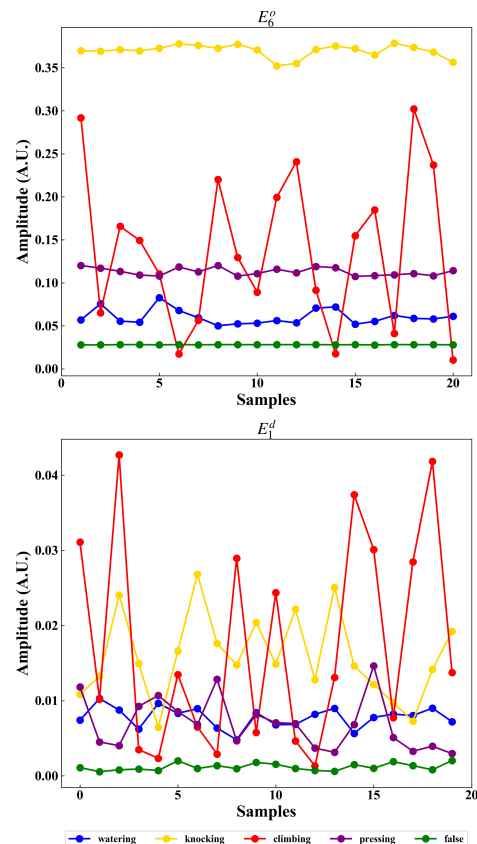


FIGURE 8. The two features for the five disturbance events.

It can be seen that our proposed method achieves the highest accuracy, and the time of feature extraction is slightly higher than wavelet energy spectrum. This is mainly because in our experiment, we use wavelet energy spectrum method to decompose the signal into two layers instead of six layers in ref [16] to obtain the highest recognition rate.

TABLE 2. The feature extraction time and average recognition accuracy of different methods.

| References | Methods | Time (s) | Accuracy (%) |
|------------|-------------------------------|----------|--------------|
| [16] | Wavelet energy spectrum | 0.06 | 73.04 |
| [33] | Time domain feature | 0.21 | 88.21 |
| [34] | Time-frequency domain feature | 0.25 | 89.96 |
| This work | EMD energy analysis | 0.11 | 93.01 |

TABLE 3. The identification results.

| Location | Acc (%) | | | | | Average |
|----------|---------|-------|-------|-------|--------|---------|
| | w | k | c | p | f | |
| 0.9km | 98.27 | 90.70 | 90.34 | 86.55 | 99.20 | 93.01 |
| 6.5km | 98.33 | 99.94 | 89.28 | 98.01 | 100.00 | 97.11 |
| 21.5km | 97.28 | 97.05 | 93.69 | 99.95 | 99.87 | 97.57 |
| Average | 97.96 | 95.90 | 91.10 | 94.84 | 99.69 | 95.90 |

(Note: w: watering, k: knocking, c: climbing, p: pressing, f: false)

B. PATTERN RECOGNITION

A total of 750 samples are obtained from five kinds of disturbance events at three locations. Then, 60% (450 samples) of the whole samples are divided randomly for training samples and the other 40% (300 samples) for testing samples. The identification process is repeated 100 times to obtain an average result, and identification accuracy (*Acc*) is shown in Table 3.

In order to further analyze performance of the classifier, the confusion matrix is obtained in Fig.9-10. On one hand, Fig.9 shows the five-category confusion matrix. The values on the diagonal are the results of correct classification. The identification rate of each disturbance signal exceeds 90%. On the other hand, the two-category confusion matrix based on alarm events (including climbing, knocking) and non-alarm events (including watering, pressing, false disturbance)

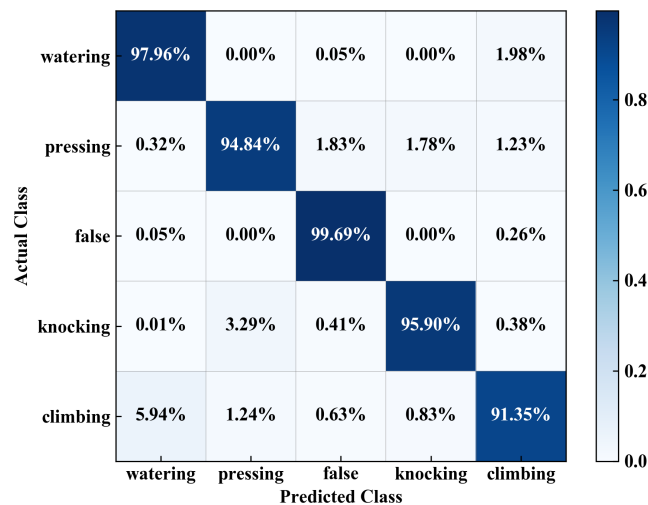


FIGURE 9. The five-category confusion matrix.

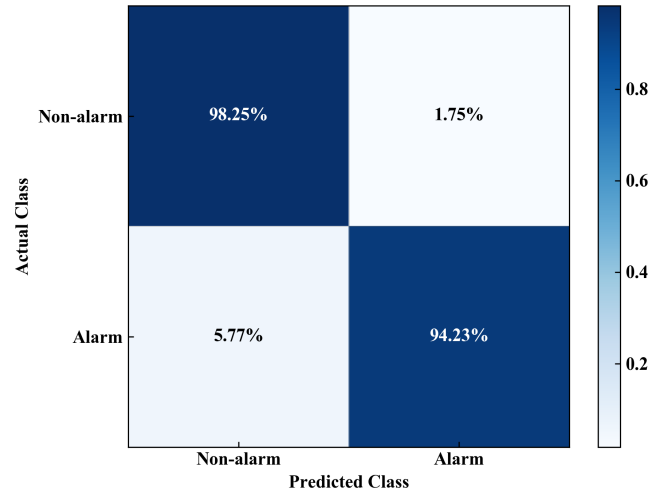


FIGURE 10. The two-category confusion matrix.

TABLE 4. The performance measures.

| Performance measure | Disturbance events | | | | |
|---------------------|--------------------|-------|-------|-------|-------|
| | w | k | c | p | f |
| <i>R</i> (%) | 97.96 | 95.90 | 91.35 | 94.84 | 99.69 |
| <i>P</i> (%) | 93.94 | 97.34 | 95.95 | 95.43 | 97.14 |
| <i>F</i> (%) | 95.91 | 96.62 | 93.59 | 95.14 | 98.40 |

(Note: w: watering, k: knocking, c: climbing, p: pressing, f: false)

are shown in Fig.10. The identification rate of the non-alarm events is as high as 98.25%, and alarm events is 94.23%.

Simultaneously, the Precision (*P*), Recall rate (*R*) and F-score (*F*) are summarized in TABLE 4. It can be seen that the *F* of watering, knocking, climbing, pressing and false disturbance is 95.91%, 96.62%, 93.59%, 95.14% and 98.40%, respectively.

We compare our proposed methods with several existing methods in the same dataset of disturbance signals, including Linear SVM, radial basis function SVM (RBF SVM), GBDT, RF and Adaboost. The recognition accuracy and recognition time of different recognition methods are shown in Table 4. Our proposed method obtains the highest value of 95.90% in terms of recognition accuracy and the recognition time of our method is 0.093 s.

The comparison of our proposed methods with other classifiers can also be illustrated by visualizing the decision boundaries. The distribution of five events at the location of 0.9 km in feature space is shown in Fig.11 (a), in which the horizontal and vertical axes represent the features E_6^o and E_1^d introduced in section IV. A respectively; five kinds of color points correspond to five different disturbance signals where blue, yellow, red, purple and green points represent ‘watering’, ‘knocking’, ‘climbing’, ‘pressing’ and ‘false’ separately. It can be seen from Fig.11(a) that watering (marked in blue) and false (marked in green) are more concentrated in feature space than other signals and thus the recognition rate of watering and false disturbance is higher than other signals. There are overlapping parts between watering and climbing (marked in red), as well as between knocking (marked in

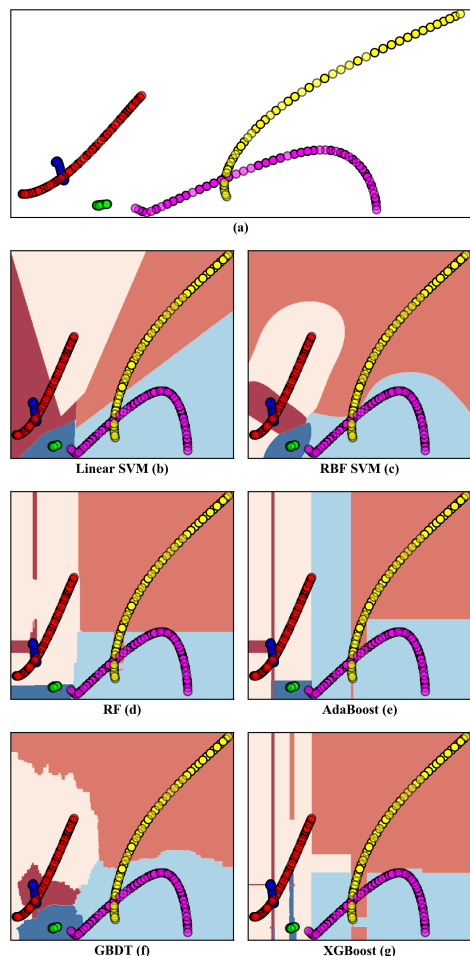


FIGURE 11. The distribution of five events (a) at the location of 0.9 km in feature space, and classification results by using visualization of decision boundaries with Linear SVM (b), RBF SVM (c), RF (d), AdaBoost (e), GBDT (f) and XGBoost (g).

yellow) and pressing (marked in purple), which explains why climbing event is identified as watering with a 5.4% probability and knocking event is identified as pressing with a 3.29% probability in Fig.9. Therefore, it can be more intuitive to display the level of recognition rates for different disturbance events by using visualization of decision boundaries.

We compare classification results with our proposed method and other classifiers including Linear SVM, RBF SVM, RF, AdaBoost and GBDT by using visualization of decision boundaries, as shown in Fig.11 (b-g), in which the different colored dots indicate that different types of signals are the same as in Fig.11 (a) and the areas with crimson, pale-yellow, brown, light blue and dark blue are separately identified by the classifier as watering, climbing, knocking, pressing and false. For the overlapping parts of ‘watering’ and ‘pressing’, Linear SVM, RBF SVM and GBDT all classify them as ‘watering’ (marked in crimson), and for the overlapping parts of ‘knocking’ and ‘pressing’, Linear SVM and GBDT classify them as ‘pressing’ (marked in light blue), and RBF SVM classifies them as ‘knocking’ (marked in brown). It can be seen from that Linear SVM, RBF SVM

TABLE 5. The recognition accuracy and recognition time of different recognition methods.

| Models | Recognition accuracy (%) | Recognition time (s) |
|------------|--------------------------|----------------------|
| Linear SVM | 88.21 | 0.102 |
| RBF SVM | 93.25 | 0.131 |
| GBDT | 94.65 | 0.322 |
| RF | 94.13 | 0.534 |
| Adaboost | 94.78 | 0.313 |
| Our method | 95.90 | 0.093 |

and GBDT cannot identify the two events with overlapping parts. Moreover, our proposed method can limit the categories of centralized distribution like ‘false’ (marked in green) to a smaller area than RF and Adaboost, which can effectively identify the category.

V. CONCLUSION

In order to reduce NAR of ϕ -OTDR, a pattern recognition method based on EMD and XGBoost is proposed in this paper. Considering non-stationarity of the signal, EMD is used for feature extraction. eleven features are extracted from normalized signals and normalized differential signals. Taking into account the impact of external environment on the system, we investigate the identification of five kinds of disturbance events including watering, knocking, climbing, pressing, false disturbance for ϕ -OTDR system with a total length of 25.05 km. The identification rates of five events (watering, knocking, climbing, pressing and false) are 97.96%, 95.90%, 91.10%, 94.84% and 99.69%, respectively. Furthermore, the classifier performance is analyzed with the help of confusion matrix. Moreover, the decision boundary visualization of our method and Linear SVM, RBF SVM, RF, AdaBoost and GBDT is obtained. Experimental results demonstrate that our proposed XGBoost method has better performance than other methods above. This method is useful to further improve the performance of ϕ -OTDR system.

REFERENCES

- [1] J. C. Juarez and H. F. Taylor, “Field test of a distributed fiber-optic intrusion sensor system for long perimeters,” *Appl. Opt.*, vol. 46, no. 11, p. 1968, Apr. 2007.
- [2] B. Dong, J. Xing, and F. Jiang, “ ϕ -OTDR optical fiber pre-warning system for perimeter security intrusion location,” *Opt. Technol.*, vol. 43, pp. 473–477, Sep. 2017.
- [3] F. Peng, N. Duan, Y.-J. Rao, and J. Li, “Real-time position and speed monitoring of trains using phase-sensitive OTDR,” *IEEE Photon. Technol. Lett.*, vol. 26, no. 20, pp. 2055–2057, Oct. 15, 2014.
- [4] W. T. Lin, S. Q. Lou, and S. Liang, “Fiber-optic distributed vibration sensor for pipeline pre-alarm,” *Appl. Mech. Mater.*, vol. 684, pp. 235–239, Oct. 2014.
- [5] F. Peng, H. Wu, X.-H. Jia, Y.-J. Rao, Z.-N. Wang, and Z.-P. Peng, “Ultra-long high-sensitivity ϕ -OTDR for high spatial resolution intrusion detection of pipelines,” *Opt. Express*, vol. 22, no. 11, pp. 13804–13810, Oct. 2014.
- [6] G. Tu, X. Zhang, Y. Zhang, F. Zhu, L. Xia, and B. Nakarmi, “The development of an ϕ -OTDR system for quantitative vibration measurement,” *IEEE Photon. Technol. Lett.*, vol. 27, no. 12, pp. 1349–1352, Jun. 2015.
- [7] Y. Wang, B. Jin, Y. Wang, D. Wang, X. Liu, and Q. Bai, “Real-Time Distributed Vibration Monitoring System Using ϕ -OTDR,” *IEEE Sensors J.*, vol. 17, no. 5, pp. 1333–1341, Dec. 2017.

- [8] Y. Muanenda, S. Faralli, C. J. Oton, and F. Di Pasquale, "Dynamic phase extraction in a modulated double-pulse ϕ -OTDR sensor using a stable homodyne demodulation in direct detection," *Opt. Express*, vol. 26, no. 2, pp. 687–701, Jan. 2018.
- [9] M. Tateda and T. Horiguchi, "Advances in optical time domain reflectometry," *J. Lightw. Technol.*, vol. 7, no. 8, pp. 1217–1224, Jun. 1989.
- [10] J. C. Juarez, E. W. Maier, K. Nam Choi, and H. F. Taylor, "Distributed fiber-optic intrusion sensor system," *J. Lightw. Technol.*, vol. 23, no. 6, pp. 2081–2087, Jun. 2005.
- [11] S. Liang, X. Sheng, S. Lou, Y. Feng, and K. Zhang, "Combination of phase-sensitive OTDR and michelson interferometer for nuisance alarm rate reducing and event identification," *IEEE Photon. J.*, vol. 8, no. 2, pp. 1–12, Apr. 2016.
- [12] H. He, L.-Y. Shao, B. Luo, Z. Li, X. Zou, Z. Zhang, W. Pan, and L. Yan, "Multiple vibrations measurement using phase-sensitive OTDR merged with mach-zehnder interferometer based on frequency division multiplexing," *Opt. Express*, vol. 24, no. 5, pp. 4842–4855, Mar. 2016.
- [13] Q. Sun, H. Feng, X. Yan, and Z. Zeng, "Recognition of a phase-sensitivity OTDR sensing system based on morphologic feature extraction," *Sensors*, vol. 15, no. 7, pp. 15179–15197, 2015.
- [14] C. Xu, J. Guan, M. Bao, J. Lu, and W. Ye, "Pattern recognition based on enhanced multifeature parameters for vibration events in -OTDR distributed optical fiber sensing system," *Microw. Opt. Technol. Lett.*, vol. 59, no. 12, pp. 3134–3141, Dec. 2017.
- [15] C. Xu, J. Guan, M. Bao, J. Lu, and W. Ye, "Pattern recognition based on time-frequency analysis and convolutional neural networks for vibrational events in ϕ -OTDR," *Opt. Eng.*, vol. 57, no. 1, Jan. 2018, Art. no. 016103.
- [16] Y. Wang, P. Wang, K. Ding, H. Li, J. Zhang, X. Liu, Q. Bai, D. Wang, and B. Jin, "Pattern recognition using relevant vector machine in optical fiber vibration sensing system," *IEEE Access*, vol. 7, pp. 5886–5895, 2019.
- [17] Y. Shi, Y. Wang, L. Zhao, and Z. Fan, "An event recognition method for ϕ -OTDR sensing system based on deep learning," *Sensors*, vol. 19, no. 15, p. 3421, Aug. 2019.
- [18] Y. Lei, J. Lin, Z. He, and M. J. Zuo, "A review on empirical mode decomposition in fault diagnosis of rotating machinery," *Mech. Syst. Signal Process.*, vol. 35, nos. 1–2, pp. 108–126, Feb. 2013.
- [19] E. Delechelle, J. Lemoine, and O. Niang, "Empirical mode decomposition: An analytical approach for sifting process," *IEEE Signal Process. Lett.*, vol. 12, no. 11, pp. 764–767, Nov. 2005.
- [20] I. Daubechies, "The wavelet transform, time-frequency localization and signal analysis," *IEEE Trans. Inf. Theory*, vol. 36, no. 5, pp. 961–1005, Oct. 1990.
- [21] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, Aug. 2016, pp. 785–794.
- [22] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Ann. Statist.*, vol. 29, no. 5, pp. 1189–1232, Oct. 2001.
- [23] L. T. Le, H. Nguyen, J. Zhou, J. Dou, and H. Moayedi, "Estimating the heating load of buildings for smart city planning using a novel artificial intelligence technique PSO-XGBoost," *Appl. Sci.*, vol. 9, no. 13, p. 2714, 2019.
- [24] N. Huang, S. Zhang, G. Cai, and D. Xu, "Power quality disturbances recognition based on a multiresolution generalized S-Transform and a PSO-improved decision tree," *Energies*, vol. 8, no. 1, pp. 549–572, 2015.
- [25] N. Huang, H. Peng, G. Cai, and J. Chen, "Power quality disturbances feature selection and recognition using optimal multi-resolution fast S-tTransform and CART algorithm," *Energies*, vol. 9, no. 11, p. 927, 2016.
- [26] J. T. Townsend, "Theoretical analysis of an alphabetic confusion matrix," *Perception Psychophys.*, vol. 9, no. 1, pp. 40–50, Jan. 1971.
- [27] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *Proc. 14th Int. Joint Conf. Artif. Intell.*, Aug. 1995, pp. 1137–1143.
- [28] H. Zhu, J. Huang, and X. Tang, "Comparing decision boundary curvature," in *Proc. 17th Int. Conf. Pattern Recognit. (ICPR)*, 2004, pp. 450–453.
- [29] E. Osuna, R. Freund, and F. Girosi, "An improved training algorithm for support vector machines," in *Proc. IEEE Signal Process. Soc. Workshop*, Feb. 1997, pp. 276–285.
- [30] Z. Wen, J. Shi, B. He, J. Chen, K. Ramamohanarao, and Q. Li, "Exploiting GPUs for efficient gradient boosting decision tree training," *IEEE Trans. Parallel Distrib. Syst.*, vol. 30, no. 12, pp. 2706–2717, Dec. 2019.
- [31] N. Huang, D. Wang, L. Lin, G. Cai, G. Huang, J. Du, and J. Zheng, "Power quality disturbances classification using rotation forest and multi-resolution fast S-transform with data compression in time domain," *IET Gener., Transmiss. Distrib.*, vol. 13, no. 22, pp. 5091–5101, Nov. 2019.
- [32] F. Wang, Z. Li, F. He, R. Wang, W. Yu, and F. Nie, "Feature learning viewpoint of AdaBoost and a new algorithm," *IEEE Access*, vol. 7, pp. 149890–149899, 2019.
- [33] X. Wang, Y. Liu, S. Liang, W. Zhang, and S. Lou, "Event identification based on random forest classifier for ϕ -OTDR fiber-optic distributed disturbance sensor," *Infr. Phys. Technol.*, vol. 97, pp. 319–325, 2019.
- [34] H. Jia, S. Liang, S. Lou, and X. Sheng, "A k-nearest neighbor algorithm-based near category support vector machine method for event identification of ϕ -OTDR," *IEEE Sensors J.*, vol. 19, no. 10, pp. 3683–3689, Jan. 2019.



ZHANDONG WANG is currently pursuing the master's degree with the School of Electronic and Information Engineering, Beijing Jiaotong University, Beijing, China.



SHUQIN LOU is currently a full-time Professor with the School of Electronic and Information Engineering, Beijing Jiaotong University, Beijing, China. She has authored or coauthored more than 200 journal articles and 60 conference papers. She holds 16 authorized patents. Her current research interests include microstructured fibers, fiber components, fiber lasers, and fiber sensors.



SHENG LIANG received the Ph.D. degree in precision instrument and machinery from Beihang University, Beijing, China, in 2011, for the work on fiber-optic distributed sensors. He is currently an Associate Professor with the Department of Physics, School of Science, Beijing Jiaotong University, Beijing. His research interests include photonics, fiber optics, and fiber-optic sensors.



XINZHI SHENG is currently a full-time Professor with the School of Science, Beijing Jiaotong University, Beijing, China. He has authored or coauthored over 80 journal articles and conference papers. He holds more than 20 authorized patents. His current research interests include optical fiber communication, fiber sensing, and fiber devices.