# Wind turbine fault detection and isolation robust against data imbalance using k-Nearest Neighbors (kNN)

| | |
|---|---|
| Journal: | *Energy Science & Engineering* |
| Manuscript ID | ESE-2023-06-0358 |
| Wiley - Manuscript type: | Original Article |
| Search Terms: | Wind energy, Power generation, Electric power systems |
| Abstract: | Due to the difficulties of system modeling, nonlinearity effects, uncertainties, and the availability of wind turbines SCADA system data, data-driven Fault Detection and Isolation (FDI) methods for Wind Turbines (WTs) have received increasing attention. In this paper, using the wind turbine SCADA data, an effective FDI scheme is proposed using the KNN classifier. The operational dataset is labeled by the status and warning datasets, and the labelled operational dataset, after eliminating invalid data, feature selection, and standardization, is used for training and validation of the FDI model. Data imbalance, which is common in real data sets, almost does not affect the accuracy of the proposed method, hence there is no need for data balancing methods in this algorithm and the accuracy is not affected by false alarms. Therefore, the proposed method has provided impressive accuracy in fault detection and isolation compared to previous research on this dataset. Also, many of the fault classes addressed in this paper were not considered in previous works on this dataset. |

SCHOLARONE™
Manuscripts

# Wind turbine fault detection and isolation robust against data imbalance using KNN

Ali Fazli[1], Javad Poshtan[1,*]

[1]Department of Electrical Engineering, Iran University of Science and Technology, Tehran, Iran
aliifazli1372@gmail.com
jposhtan@iust.ac.ir

**Corresponding Author:**
Javad Poshtan
Iran University of Science and Technology
Narmak 16846-13114, Tehran, Iran
Tel: +98(21) 73225679
Fax: +98(21) 73225777
Email: jposhtan@iust.ac.ir

**Running Title:** Wind turbine fault detection using KNN

# Wind turbine fault detection and isolation robust against data imbalance using KNN

**Abstract.** Due to the difficulties of system modeling, nonlinearity effects, uncertainties, and the availability of wind turbines SCADA system data, data-driven Fault Detection and Isolation (FDI) methods for Wind Turbines (WTs) have received increasing attention. In this paper, using the wind turbine SCADA data, an effective FDI scheme is proposed using the KNN classifier. The operational dataset is labeled by the status and warning datasets, and the labeled operational dataset, after eliminating invalid data, feature selection, and standardization, is used for training and validation of the FDI model. Data imbalance, which is common in real data sets, almost does not affect the accuracy of the proposed method, hence there is no need for data balancing methods in this algorithm and the accuracy is not affected by false alarms. Therefore, the proposed method has provided impressive accuracy in fault detection and isolation compared to previous research on this dataset. Also, many of the fault classes addressed in this paper were not considered in previous works on this dataset.

**Keywords:** Fault Detection and Isolation, KNN, SCADA data, Wind Turbine, Data imbalance

## 1. Introduction

Wind energy which is renewable and compatible with the environment is developing rapidly and in this regard, large-scale WTs are used in various countries. Therefore, in recent years, the monitoring and maintenance of WTs have attracted much attention. In a harsh environment, fault detection and condition monitoring are necessary for WTs safety and reliability [1]. Fault scenarios in WTs include sensor faults, actuator faults, and system faults [2]. Common faults and failures in WTs and the various methods that researchers use to diagnose them have been investigated and classified in [1], [3], [4]. Various wind turbine fault prognosis approaches including model-based, data-driven, knowledge-based, Artificial Intelligence (AI)-based, stochastic-based, and hybrid methods are reviewed in [4], [5]. From one perspective, fault detection and isolation methods in WTs are divided into three main categories: model-based, data-driven, and hybrid approaches. In model-based methods, a mathematical model of the WT or its subsystem produces a simulated output based on the measured input signal. By comparing the actual output of WT with the estimated output, a residual is produced to be used for fault detection and isolation [6]. Due to the system modeling difficulties and the availability of sensors data, data-driven methods in fault detection of WTs have gained increasing interest compared to the model-based approaches [7]. Data-driven FDI methods in wind turbines include acoustic emission analysis, noise analysis, oil analysis, vibration analysis, machine learning (data mining) methods, and hybrid methods [1]. Major disadvantages of acoustic emission analysis and oil analysis methods are the need for

37 background noise to be shielded and limitation to bearings with closed-loop oil supply system,
38 respectively [8]. Vibration analysis [9], [10] which is an effective tool for the condition monitoring
39 and fault diagnosis of WT drivetrains [11], includes synchronous resampling, Hilbert transform,
40 wavelet transform, statistical analysis, envelope analysis, fast Fourier transform (FFT), and short-
41 term fast Fourier transform (STFT) [6]. Among these methods, FFT is the most important tool,
42 and wavelet analysis has become increasingly popular [12]. Despite the reliability and
43 standardization of vibration analysis methods, they are expensive, intrusive, subject to sensor
44 failures, and have limited performance for low-speed rotation [8].

45 Machine learning algorithms are divided into two main categories: supervised learning which
46 consists of regression-based and classification-based approaches, and unsupervised learning
47 algorithms such as clustering [13]. If the dataset is labeled, supervised learning algorithms are
48 used, otherwise, unsupervised learning could be a good solution. If labels are categorical,
49 classifiers are used, otherwise, if labels are numerical, regression must be used [13]. In regression
50 methods, which are divided into parametric and non-parametric modeling techniques, an artificial
51 intelligence (AI) model learns the dynamics of a WT or a WT subsystem from the patterns of input
52 and output collected from a dataset [6]. By comparing the predicted behavior of this model with
53 the WT measured actual behavior, possible faults can be detected from the changes in the WT
54 behavior [14]. On the other hand, in the classification methods, an AI model is trained to recognize
55 the patterns (e.g., mode, location, and severity) of different WT faults from the input signals
56 containing the faults information [6]. The main steps of model training in machine learning are
57 data acquisition, preprocessing, feature selection and extraction, model selection, and validation;
58 specifically, these steps for classification are data acquisition, preprocessing, equalization of the
59 classes, feature selection and extraction, model fitting, cross-validation, and using the best model
60 based on validation [13]. A deep learning model is used for feature extraction in [15]. The amount
61 of normal data is much more than that of abnormal data in SCADA datasets, which makes FDI
62 models tend to be biased toward the majority class, i.e., normal data, leading to poor accuracy in
63 diagnosing faults [16]. Several methods such as setting misclassification cost, undersampling
64 majority class, oversampling minority class (e.g. synthetic minority oversampling technique
65 (SMOTE) and adaptive synthetic sampling method (ADASYN) [17]), Tomek links, and cluster
66 centers are used to overcome data imbalance challenge [18]. However, the main focus is to use an
67 appropriate method from the existing methods or develop new algorithms according to the need,
68 because if one technique works well on an imbalanced dataset, it may not work for another
69 imbalanced data set [19].

70 The main machine learning models already used in the wind turbine FDI problem are Artificial
71 Neural Networks (ANNs), Support Vector Machines (SVM), classifier fusion, and ensemble
72 classifiers. The most common types of neural networks in this area are Multi-Layer Perceptron
73 (MLP), Convolutional Neural Networks (CNN), and Recurrent Neural Networks (RNN) [20].
74 Wind turbine SCADA data includes multivariate time series that have temporal correlations within
75 each sensor variable and spatial correlations between different sensor variables. To effectively

76  capture these correlations, in [21] two parallel modules are used for feature extraction; first a Multi-
77  Scale Deep Echo State Network (MSDeepESN) for extracting temporal multi-scale features, and
78  second, a Multi-Scale Residual Network (MSResNet) that can extract spatial multi-scale features.
79  Due to the extreme data imbalance, the focal loss function is used in [21] to reduce its effect on
80  model training. In [22], first, a Multi-Kernel Fusion Convolution Neural Network (MKFCNN) is
81  designed to extract multi-scale spatial correlations of different features, then the Long Short-term
82  Memory (LSTM) is used for learning the temporal correlations among the learned spatial features.
83  After applying a down-sampling method to the training data in [22] to balance normal and faulty
84  samples, only 210 observations from 24108 normal samples are used for the training process. A
85  novel fault diagnosis framework is presented in [23] based on an end-to-end LSTM model, to learn
86  features directly from multivariate time-series data, and capture long-term dependencies through
87  the recurrent behavior and gates mechanism of LSTM. Although using temporal correlations for
88  the training and validation of the FDI system can provide useful information, it leads to the
89  problem of the algorithm's dependence on historical data for fault detection and isolation. Taking
90  into account two new fault scenarios, in [24], time-domain signals recorded from wind turbine
91  simulations are portrayed as images and fed to a multi-channel CNN, which is used for the first
92  time in wind turbine fault detection and isolation. Using a wind turbine real SCADA data based
93  on the support vector machine learning method, a fault diagnosis and prediction model is presented
94  in [18], and despite the use of various data balancing methods, the desired accuracy is not reached
95  because of the lack of attention to the fact that transient observations before and after fault
96  occurrence must be removed. As an alternative to the SVM technique, the GPC, as a Bayesian
97  nonparametric classification method, which gets away from any assumptions about the structural
98  relationship between inputs and the output, provides probabilistic information about the fault
99  types, which is valuable for making maintenance plans [25]. Using an under-sampling method,
100 just 460 samples are representative of the 24108 normal observations of the training data.
101 Therefore, a lot of information is lost by ignoring a large number of samples in the model training
102 process in [18], [22], [25]. A data-driven FDI scheme is proposed in [26] based on the fusion of
103 the decision tree, SVM, KNN, Radial Basis Function (RBF), and MLP classifiers. Although this
104 algorithm is robust against different operational conditions and measurement errors, according to
105 the confusion matrix of [26], it is biased for some minority classes toward the majority class. On
106 the other hand, using data balancing methods may prevent faulty observations from being detected
107 as normal samples, but it can result in false alarms according to confusion matrices of [18], [21],
108 [22], [25]. Figure 1 shows various methods of detection and isolation of WT faults.
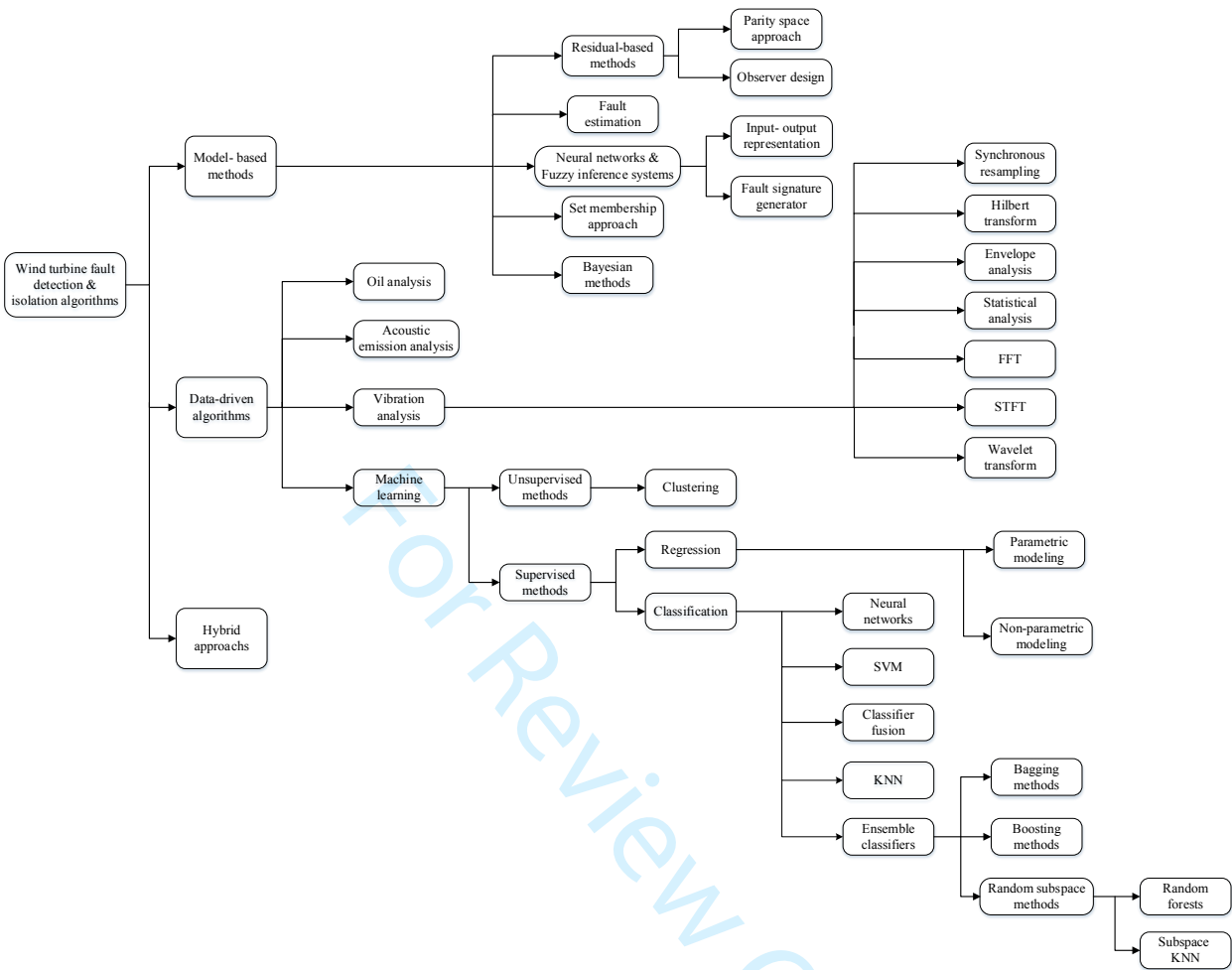
109

Figure 1. Various methods used for WT fault detection and isolation

In this paper, an FDI method is presented using wind turbine real SCADA data. First, operational SCADA data is labeled using messages of the status and warning datasets. Note that each status and warning dataset includes various faults, while in previous researches only some of the status dataset faults have been considered. Second, informative features are selected from the original features set to form an appropriate feature set. Third, the features selected are standardized in such a way that each feature has zero mean and unit standard deviation, and then the preprocessed operational dataset acts as the KNN input. Finally, WT fault detection and isolation is performed and the holdout validation method verifies remarkable accuracy of the KNN method. The motivation of this paper is to present a data-driven method for wind turbine fault detection and isolation so that the problem of data imbalance as an important challenge of the literature is solved using the potential of machine learning methods without using data balancing techniques. As an example, in the under-sampling method only some sample observations of the majority class, as a representative of this class, are used for the training process, and it is obvious that these samples can not fully represent the behavior of the no-fault (majority) class. Using data balancing methods for handling imbalanced datasets can lead to false alarms, as it can be seen in the confusion

127  matrices of [18], [21], [22], [25].  This is while the proposed algorithm can perform the training
128  and validation processes without the requirement for these interface methods and with the benefit
129  of the capacity of all dataset observations. The contributions of this paper are:

130  1)  The Model is almost robust against imbalanced data, which prevents the classifier from
131      tending to the majority class. Imbalanced classification is one of the most important
132      challenges of the wind turbine FDI methods because faults occur infrequently compared to
133      the healthy operation.
134  2)  Due to the algorithm's robustness against imbalanced data, there is no need for data
135      balancing methods. Therefore, the classifier's accuracy is not affected by false alarms.
136  3)  These two advantages make the algorithm have excellent accuracy to detect and isolate
137      wind turbine faults.
138  4)  Although the collected dataset which is used for training the FDI system is offline, the
139      algorithm determines the labels of the new data observations online.
140  5)  The proposed fault detection and isolation system determines the label of the new data
141      observation only with its current feature values, without any requirement for historical data.
142  6)  Existing SCADA data is used to train and validate the algorithm, so there is no extra cost
143      to collect and record the data.
144  7)  A large number of status and warning faults are considered (in comparison with previous
145      works on this SCADA data [18], [21], [22], [25]).
146  8)  The dependence between different faults based on physical facts and the overlap of these
147      faults' occurrence is analyzed.

148  The article structure is as follows: Section 2 describes the wind turbine SCADA dataset. In section
149  3 data preprocessing, including labeling, eliminating invalid data, feature selection, and
150  standardization is explained. Section 4 presents the parameter selection of the WT fault detection
151  and isolation model, and Section 5 includes KNN training and validation results. Finally, in section
152  6 conclusion and future work are discussed. Figure 2 shows the flowchart of the proposed scheme
153  for WT fault detection and isolation.
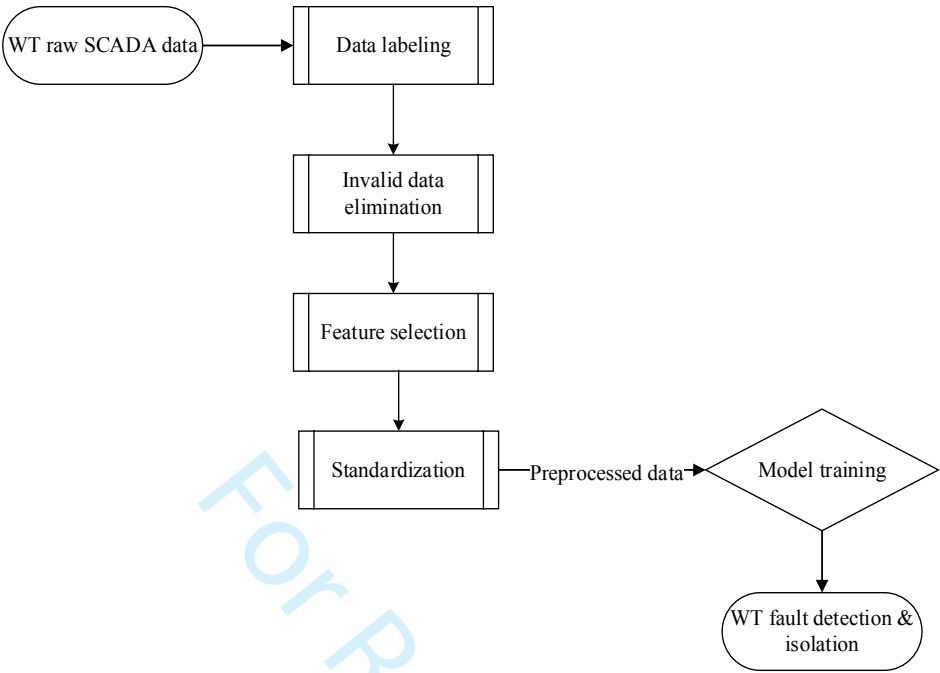
154



155

156                Figure 2. Flowchart of the proposed scheme for WT fault detection and isolation.

## 2. Data description

158    The real SCADA data in this study comes from a 3-MW direct-drive wind turbine that supplies
159    the required power of a major biomedical manufacturing plant located off the southern coast of
160    Ireland [18], [21], [22], [25], [27]. This 3-blade wind turbine comprises common major
161    components, including main shaft, generator, rotor, turbine blades, yaw system, pitch system,
162    control and power electronics system, hydraulic and cooling system, etc. The data include
163    observations over an 11-month period from May 2014 to April 2015 which contains three separate
164    datasets:

165        1) Operational data including 49027 samples with a 10-minute time interval and 66 features
166           such as ambient parameters (e.g. wind speed and ambient temperature), power
167           measurements (e.g. active power and reactive power), WT components recorded
168           temperature (e.g. tower temperature, nacelle cabinet temperature, etc.), and operating
169           conditions (e.g. nacelle position including cable twisting, operating hours).
170        2) Status data that represents the operating conditions of the wind turbine and is divided into
171           two sets of data: wind energy converter (WEC) and remote terminal unit (RTU). The WEC
172           data includes status messages that are directly related to the WT itself, while the RTU data
173           corresponds to the power control data at the point of connection to the grid, such as active
174           and reactive power. In the status dataset, a status message with a new timestamp is
175           generated each time the status changes. Therefore, it is assumed that the turbine operates

176      in that mode until the next status message is generated. Each turbine state has a "main
177      status" code and a "sub-status" code, where each main status code greater than zero
178      indicates an abnormal condition and not necessarily a fault; For example, status code 2
179      ("lack of wind") indicates the lack of sufficient wind for normal operation of a wind turbine.
180   3) Warning data that corresponds to the general information about the turbine and is not
181      directly related to the turbine operation or safety. Warning messages are about potentially
182      developing faults in the turbine. These messages have timestamps in the same way as status
183      messages and also have a "main warning" code and a "sub-warning" code. This data is also
184      divided into two sets of WEC data and RTU data. The WEC data includes warning
185      messages that are directly related to the turbine itself, while the RTU data is related to the
186      power control data at the point of connection to the grid.

## 3. Preprocessing

188 Before applying the FDI model to the operational data, as its input, some steps must be taken. First,
189 considering the messages of the status and warning datasets, the operational dataset is labeled.
190 Then some observations must be removed to prevent feeding the model with invalid data samples.
191 In the next step, non-informative features are removed and z-score normalization is applied to the
192 selected set of features. Note that one of the most important preprocessing steps in classification-
193 based FDI problems, called data balancing, is eliminated by taking advantage of using the KNN
194 model.

### 3.1. Data labeling

196 The wind turbine operational dataset is labeled based on the information in the status and
197 warning datasets so that the observations during faults occurrence are distinguished from other
198 observations. Table 1 shows status data faults along with status and sub-status codes, and the
199 symbols assigned to each fault class so that fault classification and evaluation are performed
200 easily. Faults that have not been addressed in previous works on this dataset are highlighted in
201 green. Feeding faults refer to faults in the power feeder cables of the WT, excitation errors are
202 mainly due to problems with the generator excitation system, malfunction air-cooling indicates
203 problems in the air circulation and internal temperature circulation in the WT, mains failure is
204 related to problems with mains electricity supply to the WT, and generator heating faults refer
205 to the generator overheating. Other faults are WT cable twist, inverter over-temperature, WT
206 tower transversal oscillation, fan-inverter malfunction, and yaw control fault. Among these,
207 Feeding fault and Excitation error have overlaps in 95 samples. Also, 12 samples have the
208 labels of Malfunction air cooling and Mains failure at the same time. As these overlaps cannot
209 be ignored and include a significant proportion of the number of samples of each fault, either
210 multi-label classification methods should be used or the overlap section should be considered
211 a separate class. The second solution is adopted, and for this purpose, the modified classes are
212 symbolized as EFF and AMF, respectively. In Table 2 several faults extracted from the WT
213 warning data and used to label the operational data are shown. None of these faults have been
214 investigated in previous researches on this dataset and hence all are highlighted in green.

215  According to the symbols defined for each fault, Table 3 presents the number of overlapping
216  samples between the faults defined in Tables 1 and 2. Overlaps with less than 10% of both
217  faults samples are negligible and not used in the training and testing of the proposed method.
218  For example, without considering overlap samples, the classes BF and FA with 144 and 271
219  samples, respectively, have 49 common samples, which exceeds 10% of each class sample
220  size; hence this overlap cannot be ignored. On the other hand, some faults in nature are so close
221  together that have the same main code. Accordingly, faults that are conceptually very close to
222  each other can also be considered as an integrated class. Therefore, two types of class
223  integration have been considered in this paper: first, class integration based on the number of
224  common (overlap) samples in the operational data, and second, class integration based on the
225  nature of faults. Table 4 lists the faults that are merged accordingly. Therefore, the final fault
226  labels used for data labeling are given in Table 5 and the relative frequency of these classes is
227  shown in Figure 3 using a pie chart. Figure 4 shows the scatter plot of the wind turbine's faulty
228  samples on the power curve according to labels presented in Table 5.



229
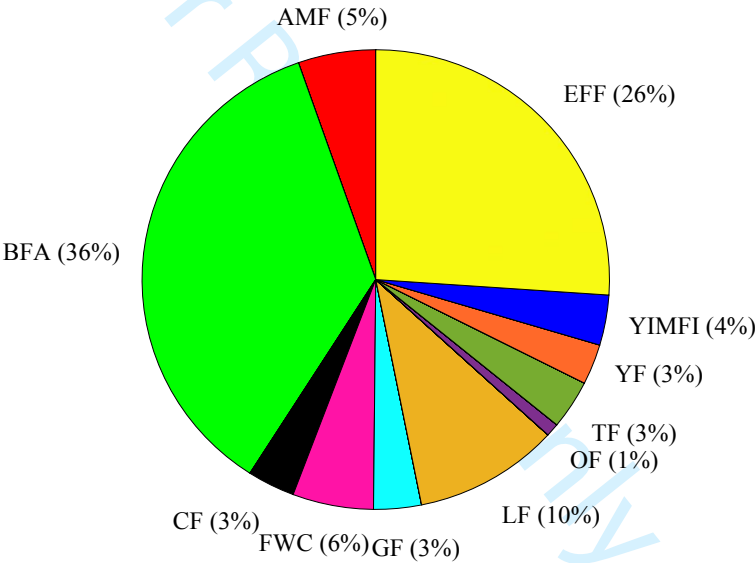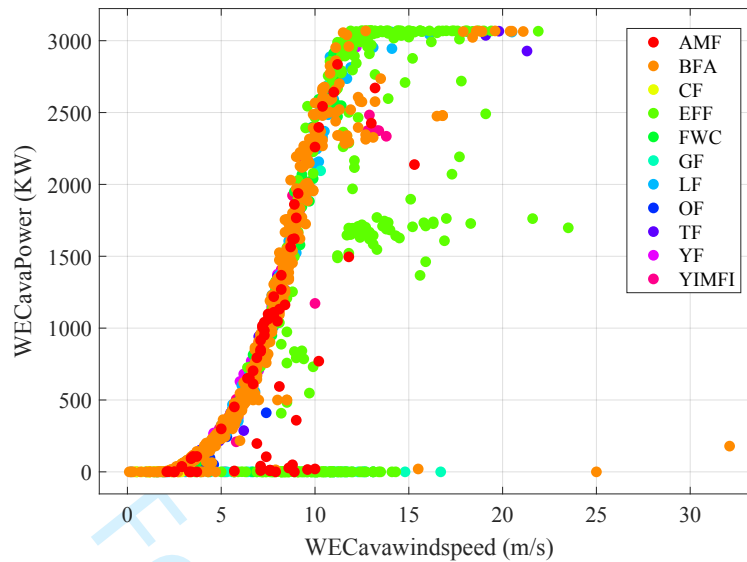230                           Figure 3. Relative frequency of the WT considered faults

Figure 4. Scatter plot of the wind turbine faulty operational data on the power curve

231

232

233    3.2.            Invalid data elimination

234    Some observations which indicate abnormal behavior of the wind turbine are not labeled with
235    fault labels, because many of these are not associated with a fault, e.g. status code 2 –"lack of
236    wind". Some of these abnormal conditions such as warning code 230 –"Power limitation
237    (10h)" are not considered as the no-fault class in data labeling. According to [22], at least 120
238    minutes before the change from the normal state and 30 minutes after the change to the normal
239    state should be considered as transient state data. This is due to the fact that, unlike in
240    simulation, in a real SCADA dataset, the change from normal to abnormal performance (and
241    vice versa) occurs with soft behavior, so the data in this transition state should not be
242    considered as normal WT performance data. In this paper with less conservative data cleaning,
243    at least 60 minutes before the change from the normal state and 20 minutes after the change to
244    the normal state are considered as transient observations. Also, sometimes a fault or a group
245    of faults occur intermittently with a small time interval, which indicates that the problem of
246    the turbine is not solved between consecutive faults. Therefore after filtering invalid normal
247    data, we have 34472 observations as the no-fault class which is labeled with the symbol NF.
248    Note that in [21], [22], [25] only 32056 observations are considered as normal class. This
249    means that in this paper, the normal class is more in line with the conditions provided by the
250    SCADA wind turbine dataset. Figure 5 shows the scatter plot of the wind turbine's normal
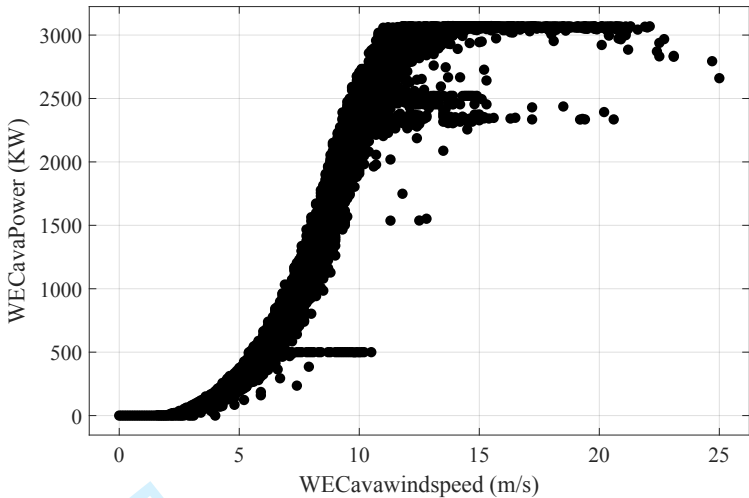251    operational data on the power curve.

252
253    Figure 5. Scatter plot of the wind turbine normal operational data on the power curve

254    ## 3.3.        Feature selection and standardization

255    Some of the 66 variables of operational data must be used in the fault detection and isolation
256    process. In short, the following steps are taken to select appropriate features, resulting in the
257    final set of features displayed in Table 6.

258    ✓ Date Time is not a numeric feature and is eliminated.
259    ✓ Variables such as Time, WEC: Operating Hours, WEC: Production kWh, and WEC:
260       Production minutes that are accumulative features result in time-based fault detection
261       and isolation if considered in the training and validation process. Therefore, even if the
262       accuracy obtained is high, it will not be reliable. For example, if a specific fault occurs
263       in March in the training set, the test sample can only be correctly predicted if it occurs
264       in the same month. Hence, for preventing the algorithm from time-based bias these
265       features are not considered in the FDI process.
266    ✓ For variables recorded as minimum, maximum, and average, it is better to use only the
267       average because the minimum and maximum values may be affected by outliers and
268       deteriorate the accuracy of the algorithm. Wind speed, Rotation, Power, and Reactive
269       power are of this type.
270    ✓ Some parameters including Sys 1 inverter cabinet temp, Bearing temp, Pitch cabinet
271       blade temp, Blade temp, Rotor temp, Stator temp, Nacelle ambient temp, and Inverter
272       cabinet temp are recorded by more than one sensor (seven, two, three, three, two, two,
273       two, and two respectively). Therefore, to prevent the curse of dimensionality, which is
274       an important challenge in the KNN algorithm, these variables are merged into one
275       variable by averaging. For example, Front bearing temp and Rear bearing temp are
276       variables that result in one variable (Bearing temp) by averaging.
277    ✓ Some non-informative features that reduce the accuracy of the algorithm are ignored.

278

279    Numerical features must be Standardized before feeding to the KNN algorithm to contribute
280    equally to the similarity measures. Standardization (z-score normalization) is a scaling method
281    where the values are centered around the mean with a unit standard deviation. This means that
282    the mean of the attribute becomes zero, and the resultant distribution has a unit standard
283    deviation, that is:

284
$$X_i^{'} = \frac{X_i - \mu_i}{\sigma_i} \tag{2}$$

285    where $X_i$ represents the $i$-th original feature, $\mu_i$ and $\sigma_i$ are its mean and standard deviation
286    respectively, and $X_i^{'}$ is the standardized $i$-th feature for $i=1,\ldots,25$.

287

288    ## 4.  KNN parameter selection

289    Considering the final fault classes, presented in Table 5, the preprocessed operational data is used
290    to determine the FDI model parameter based on the classification error. A KNN classifier is trained
291    for different values of the parameter K. Figure 6 shows the effect of changes in this parameter on
292    the classification error by the 10-fold validation method. According to this figure, increasing the
293    value of parameter K results in increasing the classification error, so $K = 1$ is the best choice for
294    this classifier. Choosing this value usually can yield overfitting, hence to avoid this, the holdout
295    validation is adopted in the procedure of the wind turbine FDI model design. As shown in Figure
296    6, the KNN classifier can well classify different classes of wind turbine's imbalanced SCADA
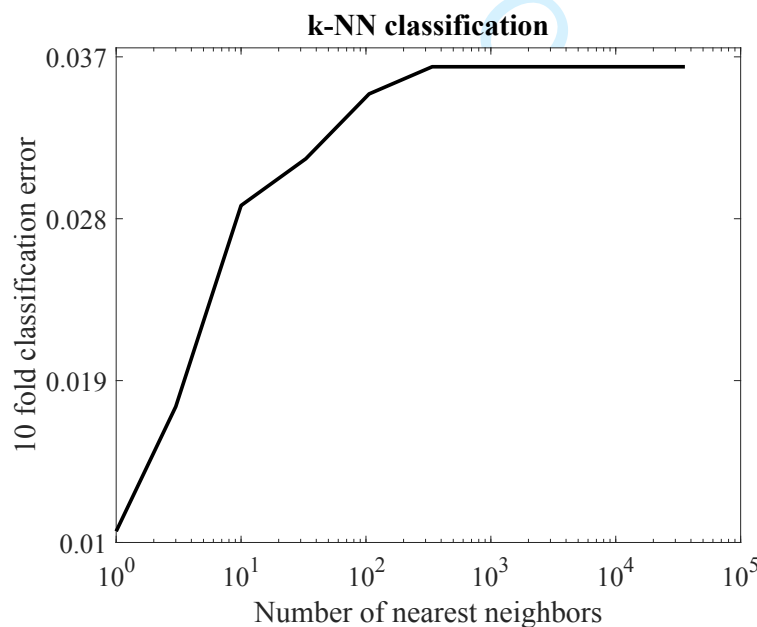297    dataset.

298


299    Figure 6. Classification error for different values of parameter *K*

300 Many references confirm that the performance of KNN is affected by data imbalance [28]–[30].
301 Conceptually the KNN algorithm calculates the (Euclidian) distances of the training dataset
302 observations from a validation sample (that is to be labeled) and assigns the majority label of $K$
303 nearest neighbors of the validation sample to it. Therefore, if the observations of each class (in the
304 feature space of the dataset) have little distance from each other and a small parameter $K$ is
305 selected, the performance of this algorithm will not be affected by data imbalance. In other words,
306 if the distance of samples within a class is less than the distance between samples of different
307 classes, and a small parameter $K$ is selected, this algorithm will be robust against data imbalance.
308 For example, consider that one class of the dataset has 1000 samples and another class has only 2
309 observations, but these 2 samples are well close to each other in the feature space. In this case,
310 considering $K = 1$, the KNN algorithm will assign a true label to those two samples and will not
311 be affected by the data imbalance. Also, according to Figure 6, by increasing the parameter $K$, the
312 error reaches 0.03645, which means that the algorithm tends to the majority class: Note that
313 0.03645 * 35776=1304, exactly equivalent to the false prediction of the minority classes samples
314 (1304 fault observations).

## 5. Training and validation

316 After choosing suitable parameters according to the previous section's explanations, the KNN
317 model is applied to the train data, which is 80% of the data. Therefore, 20% of the data is used as
318 validation data, employing the holdout validation method, thus none of the validation data samples
319 are involved in the training process. Considering the faults given in Table 5, the KNN model results
320 in an impressive performance as shown in Figure 7. This confusion chart illustrates that almost all
321 samples of the validation data are assigned properly to the true classes. The following parameters
322 are used to clarify this performance:

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} \tag{3}$$

$$TPR_{ave} = \frac{\sum_{m=1}^{C} TPR_m}{C} \tag{4}$$

$$PPV_{ave} = \frac{\sum_{m=1}^{C} PPV_m}{C} \tag{6}$$

326 where $C$ is the number of classes including no-fault and faulty categories, $TP$ is the number of true
327 predictions of the class, $FN$ represents the number of false predictions of the class, $P$ indicates the
328 number of all samples of the class, and FP is the representative of the false alarms of the class.
329 Note that average $TPR$ in (4) and average $PPV$ in (6), which are proposed as validation parameters
330 in this paper, can be useful for the validation of the imbalance classification. The values of $TPR$
331 and $PPV$ obtained for all classes are greater than 77.8% and 75% respectively according to Figure
332 7, confirming the remarkable performance of the proposed model in the WT fault detection and

333   isolation is confirmed. Note that the training and validation data are chosen randomly based on the
334   80-20 proportions. So, each time the training and validation process takes place, the results may
335   be slightly different. In multiple runs, however, the number of false predictions does not exceed
336   75, which is a good achievement compared to the previously reported methods. Figure 7 also
337   confirms this statement because almost 0.01062 10-fold classification error is equal to
338   0.01062*0.2*35776≈76 false predictions in the holdout validation. Table 7 compares the proposed
339   FDI model with related previous works on this dataset, in terms of accuracy, the number of the
340   fault classes considered, the number of the normal class samples, overall accuracy (the ratio of the
341   number of the whole true predictions to the number of the whole false predictions), average *TPR*,
342   average *PPV*, and requirement of the method to data balancing. According to the comparison made
343   in [21] and [22], different algorithms including several deep learning methods such as CNN,
344   LSTM, MKFCNN, CNN-LSTM, CNN-GRU, MSDeepESN, and MSResNet neural networks, as
345   well as shallow machine learning techniques such as random forests (RF), SVM, decision tree
346   (DT), and GPC are adopted for comparison on the present SCADA dataset. Note that overall
347   accuracy is not a good validation criterion for imbalance classification, because a case such as
348   "good *TPR* for no-fault class and bad prediction for other classes" also can result in high overall
349   accuracy.

**True Class** (rows) vs **Predicted Class** (columns)

| True \ Pred | AMF | BFA | CF | EFF | FWC | GF | LF | NF | OF | TF | YF | YIMFI |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AMF | 11 | 1 | | | | | | 2 | | | | |
| BFA | | 75 | | 2 | | | | 15 | | | | |
| CF | | | 8 | | | | | | | | | |
| EFF | | | | 65 | | | | 3 | | | | |
| FWC | | | | | 13 | | 1 | | | | | |
| GF | | | | | | 8 | | | | | | |
| LF | | | | | | | 21 | 4 | | 1 | | |
| NF | 1 | 8 | | 3 | 2 | | 2 | 6875 | | 1 | 2 | |
| OF | | | | | | | | | 2 | | | |
| TF | | | | | | | | 2 | | 7 | | |
| YF | | | | | | | | 1 | | | 6 | |
| YIMFI | | | | | | | | 1 | | | | 8 |

350
351 Figure 7. The confusion chart of the proposed model validation

352

## 6. Conclusion

354 In this paper, the KNN model is used for WT fault detection and isolation. For the purpose of
355 preprocessing, four steps of data labeling, invalid data elimination, feature selection, and
356 standardization are applied to the raw SCADA dataset of a wind turbine, respectively. It was shown
357 that the proposed method is robust against the data imbalance challenge, which prevents the

358    classifier from tending to the majority class. Due to the algorithm's robustness against imbalanced
359    data, there is no need for data balancing methods. Therefore, the classifier's accuracy is not
360    affected by false alarms. As a result of the above, the proposed model shows higher accuracy in
361    the WT fault detection and isolation, compared to previous researches on this dataset. The
362    proposed FDI system can be used for online monitoring, unlike the training step which is an offline
363    process. This system determines the label of the new data observation only with its current feature
364    values, without any requirement for historical data. Although using temporal correlations for the
365    training and validation of the FDI system can provide useful information, it leads to the problem
366    of the algorithm's dependence on historical data for fault detection and isolation. Since the existing
367    SCADA data was used to train and validate the algorithm, there is no extra cost to collect and
368    record the data. It is worth mentioning that many faults are considered, especially from the warning
369    data, which are not indicated in the previous related works. Finally, the dependence of the faults
370    is analyzed based on the nature of the faults and the overlap of their occurrence. As a future work,
371    predictive maintenance can be pursued by developing the existing model to predict wind turbine
372    faults.

## References

373

374    [1]    A. Purarjomandlangrudi, G. Nourbakhsh, M. Esmalifalak, and A. Tan, "Fault detection in wind
375           Turbine: A systematic literature review," *Wind Eng.*, vol. 37, no. 5, pp. 535–546, 2013, doi:
376           10.1260/0309-524X.37.5.535.

377    [2]    D. Yu, Z. M. Chen, K. S. Xiahou, M. S. Li, T. Y. Ji, and Q. H. Wu, "A radically data-driven method
378           for fault detection and diagnosis in wind turbines," *Int. J. Electr. Power Energy Syst.*, vol. 99, pp.
379           577–584, 2018, doi: 10.1016/j.ijepes.2018.01.009.

380    [3]    Y. Amirat, M. E. H. Benbouzid, E. Al-Ahmar, B. Bensaker, and S. Turri, "A brief status on condition
381           monitoring and fault diagnosis in wind energy conversion systems," *Renew. Sustain. Energy Rev.*,
382           vol. 13, no. 9, pp. 2629–2636, 2009, doi: 10.1016/j.rser.2009.06.031.

383    [4]    M. Rezamand, M. Kordestani, R. Carriveau, D. S. K. Ting, M. E. Orchard, and M. Saif, "Critical
384           Wind Turbine Components Prognostics: A Comprehensive Review," *IEEE Trans. Instrum. Meas.*,
385           vol. 69, no. 12, pp. 9306–9328, 2020, doi: 10.1109/TIM.2020.3030165.

386    [5]    M. Kordestani, M. Saif, M. E. Orchard, R. Razavi-Far, and K. Khorasani, "Failure Prognosis and
387           Applications - A Survey of Recent Literature," *IEEE Trans. Reliab.*, vol. 70, no. 2, pp. 728–748,
388           2021, doi: 10.1109/TR.2019.2930195.

389    [6]    W. Qiao and D. Lu, "A Survey on Wind Turbine Condition Monitoring and Fault Diagnosis - Part
390           II: Signals and Signal Processing Methods," *IEEE Trans. Ind. Electron.*, vol. 62, no. 10, pp. 6546–
391           6557, 2015, doi: 10.1109/TIE.2015.2422394.

392    [7]    G. Jiang, P. Xie, H. He, and J. Yan, "Wind Turbine Fault Detection Using a Denoising Autoencoder
393           with Temporal Information," *IEEE/ASME Trans. Mechatronics*, vol. 23, no. 1, pp. 89–100, 2018,
394           doi: 10.1109/TMECH.2017.2759301.

395    [8]    S. M. M. Aval and A. Ahadi, "Wind turbine fault diagnosis techniques and related algorithms," *Int.
396           J. Renew. Energy Res.*, vol. 6, no. 1, pp. 80–89, 2016.

397    [9]    M. He and D. He, "A new hybrid deep signal processing approach for bearing fault diagnosis using

398        vibration signals," *Neurocomputing*, vol. 396, pp. 542–555, 2020, doi:
399        10.1016/j.neucom.2018.12.088.

400    [10]   C. Liu, G. Cheng, B. Liu, and X. Chen, "Bearing fault diagnosis method with unknown variable
401        speed based on multi-curve extraction and selection," *Meas. J. Int. Meas. Confed.*, vol. 153, 2020,
402        doi: 10.1016/j.measurement.2019.107437.

403    [11]   W. Teng, X. Ding, S. Tang, J. Xu, B. Shi, and Y. Liu, "Vibration analysis for fault detection of wind
404        turbine drivetrains—a comprehensive investigation," *Sensors*, vol. 21, no. 5, pp. 1–33, 2021, doi:
405        10.3390/s21051686.

406    [12]   D. Strömbergsson, P. Marklund, K. Berglund, and P. E. Larsson, "Bearing monitoring in the wind
407        turbine drivetrain: A comparative study of the FFT and wavelet transforms," *Wind Energy*, vol. 23,
408        no. 6, pp. 1381–1393, 2020, doi: 10.1002/we.2491.

409    [13]   A. Stetco *et al.*, "Machine learning methods for wind turbine condition monitoring: A review,"
410        *Renew. Energy*, vol. 133, pp. 620–635, 2019, doi: 10.1016/j.renene.2018.10.047.

411    [14]   A. Zaher, S. D. J. McArthur, D. G. Infield, and Y. Patel, "Online Wind Turbine Fault Detection
412        through Automated SCADA Data Analysis," *Wind ENERGY*, vol. 12, no. 6, pp. 574–593, 2009, doi:
413        10.1002/we.319.

414    [15]   Y. Liu, H. Cheng, X. Kong, Q. Wang, and H. Cui, "Intelligent wind turbine blade icing detection
415        using supervisory control and data acquisition data and ensemble deep learning," *Energy Sci. Eng.*,
416        vol. 7, no. 6, pp. 2633–2645, 2019, doi: 10.1002/ese3.449.

417    [16]   L. Chen, G. Xu, Q. Zhang, and X. Zhang, "Learning deep representation of imbalanced SCADA
418        data for fault detection of wind turbines," *Meas. J. Int. Meas. Confed.*, vol. 139, pp. 370–379, 2019,
419        doi: 10.1016/j.measurement.2019.03.029.

420    [17]   A. Burkov, *The Hundred-Page Machine Learning Book*. 2019. doi:
421        10.1080/15228053.2020.1766224.

422    [18]   K. Leahy, R. L. Hu, I. C. Konstantakopoulos, C. J. Spanos, A. M. Agogino, and D. T. J. O'Sullivan,
423        "Diagnosing and Predicting Wind Turbine Faults from SCADA Data Using Support Vector
424        Machines," *Int. J. Progn. Heal. Manag.*, vol. 9, no. 1, pp. 1–11, 2018.

425    [19]   N. Rout, D. Mishra, and M. K. Mallick, "Handling Imbalanced Data : A Survey," in *International
426        Proceedings on Advances in Soft Computing, Intelligent Systems and Applications*, 2018, pp. 431–
427        443.

428    [20]   G. Helbing and M. Ritter, "Deep Learning for fault detection in wind turbines," *Renew. Sustain.
429        Energy Rev.*, vol. 98, pp. 189–198, 2018, doi: 10.1016/j.rser.2018.09.012.

430    [21]   Q. He, Y. Pang, G. Jiang, and P. Xie, "A spatio-temporal multiscale neural network approach for
431        wind turbine fault diagnosis with imbalanced SCADA data," *IEEE Trans. Ind. Informatics*, pp. 1–
432        9, 2020, doi: 10.1109/TII.2020.3041114.

433    [22]   Y. Pang, Q. He, G. Jiang, and P. Xie, "Spatio-temporal fusion neural network for multi-class fault
434        diagnosis of wind turbines based on SCADA data," *Renew. Energy*, vol. 161, pp. 510–524, 2020,
435        doi: 10.1016/j.renene.2020.06.154.

436    [23]   J. Lei, C. Liu, and D. Jiang, "Fault diagnosis of wind turbine based on Long Short-term memory
437        networks," *Renew. Energy*, vol. 133, pp. 422–432, 2019, doi: 10.1016/j.renene.2018.10.031.

438    [24]   S. Zare and M. Ayati, "Simultaneous fault diagnosis of wind turbine using multichannel

439     convolutional neural networks," *ISA Trans.*, vol. 108, pp. 230–239, 2021, doi:
440     10.1016/j.isatra.2020.08.021.

441  [25]  Y. Li, S. Liu, and L. Shu, "Wind turbine fault diagnosis based on Gaussian process classifiers
442     applied to operational data," *Renew. Energy*, vol. 134, pp. 357–366, 2019, doi:
443     10.1016/j.renene.2018.10.088.

444  [26]  V. Pashazadeh, F. R. Salmasi, and B. N. Araabi, "Data driven sensor and actuator fault detection
445     and isolation in wind turbine using classifier fusion," *Renew. Energy*, vol. 116, pp. 99–106, 2018,
446     doi: 10.1016/j.renene.2017.03.051.

447  [27]  R. L. Hu, K. Leahy, I. C. Konstantakopoulos, D. M. Auslander, C. J. Spanos, and A. M. Agogino,
448     "Using domain knowledge features for wind turbine diagnostics," in *Proceedings - 2016 15th IEEE
449     International Conference on Machine Learning and Applications, ICMLA 2016*, 2017, pp. 300–305.
450     doi: 10.1109/ICMLA.2016.172.

451  [28]  M. E. Kadir, P. S. Akash, S. Sharmin, A. A. Ali, and M. Shoyaib, "A Proximity Weighted Evidential
452     k Nearest Neighbor Classifier for Imbalanced Data," in *Pacific-Asia Conference on Knowledge
453     Discovery and Data Mining*, 2020, pp. 71–83. doi: 10.1007/978-3-030-47436-2_6.

454  [29]  K. Kim, "Normalized class coherence change-based kNN for classification of imbalanced data,"
455     *Pattern Recognit.*, vol. 120, 2021, doi: 10.1016/j.patcog.2021.108126.

456  [30]  S. Zhang, "Cost-sensitive KNN classification," *Neurocomputing*, vol. 391, pp. 234–242, 2020, doi:
457     10.1016/j.neucom.2018.11.101.

458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482                              Table 1. Status data faults

| Fault | Status code | Sub-status code | Number of samples | Symbol |
|---|---|---|---|---|
| Generator heating: Hygrostat inverter | 9 | 3 | 43 | GF |
| Feeding fault | 62 | various | 260 | FF |
| Excitation error : Overvoltage DC-link | 80 | 21 | 175 | EF |
| Timeout warn message : Malfunction air cooling | 228 | 100 | 62 | AF |
| Mains failure : Under voltage L1/ Start delay | 60 | 11/2 | 20 | MF |
| Cable twisted : Left/Right (2-3 turns) | 21 | 1/2 | 43 | CF |
| Over-temperature : Power choke inverter 4 | 67 | 604 | 12 | OF |
| Tower oscillation : Transversal oscillation | 31 | 1 | 3 | TF |
| Malfunction fan-inverter: Other control board errors system 3 | 26 | 373 | 3 | IF |
| Yaw control fault: TWK NOCN: fault in cam switch mechanism | 22 | 226 | 5 | YF |

483                                Table 2. Warning data faults

| Fault | Warning code | Sub-warning code | Number of samples | Symbol |
|---|---|---|---|---|
| Faulty yaw inverter: Rotational speed error motor x/DC-link voltage instabil system y/Feedback mains contactor faulty(-) | 25 | various | 28 | YIF |
| Yaw control fault : TWK NOCN installed | 22 | 200 | 32 | YCF |
| Fault blade load control: Blade load curve A/B/C too low (28) | 49 | 193/293/393 | 195 | BF |
| Fault air cooling: Generator pressure to low (7) | 100 | 11 | 327 | FA |
| Fault water cooling: Pressure too low (7)/ No flow rate (7) | 101 | 1/5 | 85 | FWC |
| Fault lubrication system: Grease reservoir empty (rotor) (90) | 58 | 1 | 140 | LF |
| Tower oscillation: Longitudinal oscill. (low wind speed) | 31 | 32 | 42 | TLF |
| Malfunction fan-inverter: DC-link voltage instable system 3/ Circuit breaker nacelle fan tripped | 26 | 316, 331, 335 | 26 | MFI |

484

485                      Table 3. Number of the faults overlapping samples

| Overlap fault | Number of each fault samples without considering the overlap | Number of overlap samples | Symbol | Can be ignored |
|---|---|---|---|---|
| BF+FA | 144+271 | 49 | BFA | No |
| YIF+MFI | 19+19 | 5 | YIMFI | No |
| FA+FWC | 271+73 | 5 | FAWC | Yes |
| LF+FWC | 131+73 | 7 | LFWC | Yes |
| LF+FA | 131+271 | 2 | LFA | Yes |
| BF+YIF | 144+19 | 2 | BYIF | Yes |

| | | | | |
|---|---|---|---|---|
| AF+MF | 50+8 | 12 | AMF | No |
| EF+FF | 80+165 | 95 | EFF | No |

486

Table 4. Integrated classes based on the fault nature

| Fault 1 | Number of samples | Fault 2 | Number of samples | Integrated class symbol | Number of samples |
|---|---|---|---|---|---|
| YF | 5 | YCF | 32 | YF | 37 |
| TF | 3 | TLF | 42 | TF | 45 |
| IF | 3 | MFI | 19 | MFI | 22 |

487

Table 5. Final fault labels

| Fault | Decomposition | Number of samples |
|---|---|---|
| BFA | BF+FA | 464 |
| YIMFI | YIF+MFI+IF | 46 |
| EFF | EF+FF | 340 |
| AMF | AF+MF | 70 |
| TF | TF+TLF | 45 |
| YF | YF+YCF | 37 |
| GF | GF | 43 |
| CF | CF | 43 |
| OF | OF | 12 |
| FWC | FWC | 73 |
| LF | LF | 131 |

488

489

Table 6. Final SCADA variables

| No. | Variable Name | No. | Variable Name |
|---|---|---|---|
| 1 | Ava wind speed | 14 | Nacelle ambient temp (Averaged) |
| 2 | Ava Rotation | 15 | Nacelle temp |
| 3 | Ava Power | 16 | Nacelle cabinet temp |
| 4 | Ava Nacelle's position including cable twisting | 17 | Main carrier temp |
| 5 | Ava Reactive Power | 18 | Rectifier cabinet temp |
| 6 | Ava blade angle A | 19 | Inverter cabinet temp (Averaged) |
| 7 | Sys1 inverter cabinet temp (Averaged) | 20 | Ambient temp |
| 8 | Spinner temp | 21 | Tower temp |
| 9 | Bearing temp (Averaged) | 22 | Control cabinet temp |
| 10 | Pitch cabinet blade temp (Averaged) | 23 | Transformer temp |
| 11 | Blade temp (Averaged) | 24 | Inverter averages |
| 12 | Rotor temp (Averaged) | 25 | Inverter std dev |

| 13 | Stator temp (Averaged) | 26 | Status (Label) |
|----|------------------------|-----|----------------|

490

491          Table 7. Comparison of the proposed FDI model with related previous works on this dataset

| FDI model | Number of fault classes | Number of normal class samples | Overall accuracy | Average TPR | Average PPV | The requirement for data balancing |
|-----------|-------------------------|--------------------------------|------------------|-------------|-------------|------------------------------------|
| Proposed model | 11 | 35870 | 99.27% | 90.13% | 91.71% | No |
| STMNN [21] | 5 | 32056 | 94.93% | 90.64% | 82.68% | Yes |
| STFNN [22] | 5 | 32056 | 93.84% | 84.28% | 63.74% | Yes |
| [25] | 5 | 32056 | 91.79% | 78.87% | 44.40% | Yes |
| [18] | 4 | 48503 | 88.67% | 74.99% | 40.98% | Yes |
| RF [22] | 5 | 32056 | 84.75% | 60.95% | 46.09% | Yes |
| SVM [22] | 5 | 32056 | 72.05% | 64.19% | 37.24% | Yes |
| DT [22] | 5 | 32056 | 80.54% | 62.24% | 36.57% | Yes |
| GPC [22] | 5 | 32056 | 91.80% | 79.42% | 41.86% | Yes |
| CNN [22] | 5 | 32056 | 86.90% | 62.86% | 32.94% | Yes |
| LSTM [22] | 5 | 32056 | 87.59% | 53.61% | 30.03% | Yes |
| MKFCNN [22] | 5 | 32056 | 91.11% | 57.42% | 40.80% | Yes |
| CNN-LSTM [22] | 5 | 32056 | 89.32% | 60.84% | 39.94% | Yes |
| CNN-GRU [22] | 5 | 32056 | 84.08% | 59.15% | 34.54% | Yes |
| MSDeepESN [21] | 5 | 32056 | 87.77% | 64.11% | 37.45% | Yes |
| LSTM [21] | 5 | 32056 | 78.22% | 61.87% | 33.03% | Yes |
| MSResNet [21] | 5 | 32056 | 91.49% | 65.51% | 50.69% | Yes |
| CNN [21] | 5 | 32056 | 82.04% | 56.98% | 28.50% | Yes |

492