

Northumbria Research Link

Citation: Morrison, Rory, Liu, Xiaolei and Lin, Zi (2022) Anomaly detection in wind turbine SCADA data for power curve cleaning. *Renewable Energy*, 184. pp. 473-486. ISSN 0960-1481

Published by: Elsevier

URL: <https://doi.org/10.1016/j.renene.2021.11.118>
<<https://doi.org/10.1016/j.renene.2021.11.118>>

This version was downloaded from Northumbria Research Link:
<https://nrl.northumbria.ac.uk/id/eprint/51098/>

Northumbria University has developed Northumbria Research Link (NRL) to enable users to access the University's research output. Copyright © and moral rights for items on NRL are retained by the individual author(s) and/or other copyright owners. Single copies of full items can be reproduced, displayed or performed, and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided the authors, title and full bibliographic details are given, as well as a hyperlink and/or URL to the original metadata page. The content must not be changed in any way. Full items must not be sold commercially in any format or medium without formal permission of the copyright holder. The full policy is available online: <http://nrl.northumbria.ac.uk/policies.html>

This document may differ from the final, published version of the research and has been made available online in accordance with publisher policies. To read and/or cite from the published version of the research, please visit the publisher's website (a subscription may be required.)



**Northumbria
University**
NEWCASTLE



UniversityLibrary

1 Anomaly Detection in Wind Turbine SCADA
2 Data for Power Curve Cleaning

3 Rory Morrison¹, Xiaolei Liu^{1,*}, and Zi Lin²

4 ¹*James Watt School of Engineering, University of Glasgow, Glasgow, G12 8QQ.*

5 ²*Department of Mechanical and Construction Engineering, Northumbria University,*
6 *Newcastle upon Tyne, NE1 8ST, UK.*

7 ** Corresponding author, xiaolei.liu@glasgow.ac.uk*

8 **Abstract**

9 Wind turbine power curve cleaning, by way of removing curtailment,
10 stoppage, and other anomalies, is an essential step in making raw data
11 useable for further analysis, such as determining turbine performance, site
12 characteristics, or improving forecasting models. Typically, data comes as
13 SCADA (Supervisory Control and Data Acquisition) data, so contains not
14 only environmental and turbine performance data but also the control ac-
15 tion imposed on the turbine by the operator. Many different anomaly
16 detection (AD) methods have been proposed to clean power curves; how-
17 ever, few papers have explored filtering explicit and obvious anomalies
18 from the SCADA prior to running AD. This paper actively explores this
19 filtering impact by comparing the performances of 4 different AD methods
20 with/without filtering. These are: iForest, Local Outlier Factor, Gaussian
21 Mixture Models, and k-Nearest Neighbours. Each approach is evaluated
22 in terms of prediction error, data removal rates, and ability to maintain

23 the underlying wind statistical characteristics. The results show the effec-
24 tiveness of filtering with every technique showing improvement compared
25 to its unfiltered counterpart. Furthermore, Gaussian Mixture Models are
26 shown to provide favourable accuracy whilst maintaining wind variability,
27 however, with the wide range of performances of methods, a user's choice
28 may be different depending on their needs.

29 **Keywords**— Wind turbine, power curve, data cleaning, anomaly detection

30 **Nomenclature**

31 **Abbreviations**

32 *AD* Anomaly Detection

33 *BIC* Bayesian Information Criterion

34 *CPU* Central Processing Unit

35 *DBSCAN* Density Based Spatial Clustering of Applications with Noise

36 *FIML* Full Information Maximum Likelihood

37 *GMM* Gaussian Mixture Modelling

38 *iForest* Isolation Forest

39 *IQR* Interquartile Range

40 *kNN* k Nearest Neighbours

41 *LOF* Local Outlier Factor

42 *MAR* Missing At Random

43 *MCAR* Missing Completely At Random

44 *MNAR* Missing Not At Random

45 *NN* Neural Network

46 *RMSE* Root Mean Squared Error

47 *SCADA* Supervisory Control and Data Acquisition

48 *WT* Wind Turbine

49 *WTFC* Wind Turbine Power Curve

50 **Symbols - Isolation Forest.**

51 $c(n)$ Average of $h(x)$ for n instances

52 $E(h(x))$ Average path length across all iTress

53 $h(x)$ Averaged path length

54 n Number of instances

55 $s(x, n)$ Anomaly score

56 **Symbols - Gaussian Mixture Models**

57 μ_p Mean of a given variable

58 *IQR* Interquartile range, $Q_{75} - Q_{25}$

59 k Number of mixtures assumed

60 p Number of variables

61 p_{lower}, p_{upper} Lower and upper bounds of the box plot

62 Q_{25}, Q_{75} Lower and upper quartiles, equivalent to 25th and 75th quartiles.

63 **Symbols - k Nearest Neighbours**

64 k Number of nearest neighbours

65 **Symbols - Local Outlier Factor**

66 *lrd* local reach distance

67 N_{MinPts} Number of nearest neighbours to consider

68 o A nearest neighbour of p when considering *MinPts* of nearest neighbours

69 *reach – dist* Reach distance

70 *x* The instance being studied

71 **Symbols**

72 δIQR_u Percentage difference in wind speed IQR

73 γ Elimination rate

74 e_m Prediction error as percentage of P_r

75 $h(u_i)$ Predicted power for instance i with windspeed u

76 $IQR_{u,b}, IQR_{u,a}$ IQR of wind speed before and after AD

77 n Number of instances in test set

78 N_b, N_a Number of instances before and after AD

79 p_i Actual power value for instance i

80 P_r Rated power of the wind turbine

81 u Wind speed

82 **Units**

83 *GW* Gigawatts

84 *km* Kilometer

85 *kW* Kilowatts

86 *m/s* Meters/second

87 *MW* Megawatts

88 *RPM* Revolutions per Minute

89 1 Introduction

90 The European Union and United Kingdom are committed to extensive targets to in-
91 crease offshore wind energy capacity as part of the greening of the energy sector. This
92 is to support the commitments of many nations to the Paris Climate Agreement. As
93 of November 2020, the European Union has committed to increasing their 12GW of
94 capacity to 60GW by 2030 and 300GW by 2050 [1]. Similarly the United Kingdom
95 pledged to increase their 10GW capacity to 40GW by 2030 [2]. It is certainly an
96 exciting time to be involved in wind energy academia as these targets will need to
97 be supported by research to overcome the plethora of challenges facing such ambi-
98 tious targets, such as continuing the reduction of levelized cost of energy, turbine life
99 extension, increasing reliability, and improving forecasting models, to name a few.

100 The common thread to these research topics is their reliance, in part, on *SCADA*
101 data from already deployed wind turbines. *SCADA* (Supervisory Control and Data
102 Acquisition) data is continuously generated by each wind turbine (WT) when deployed.
103 It documents production (power output), turbine parameters (rotor RPM, blade pitch
104 angle, braking, bearing temperatures etc), supervisory action imposed by the operator
105 on the WT, and environmental conditions (air temperature, wind speed, ice indication
106 etc). If we wish to use *SCADA* data to explore the relationship WTs have with the
107 environment then anomalies must be removed and the *power curve* cleaned. The
108 wind turbine power curve (WTPC) is simply a plot of wind speeds (u) against power
109 produced by the WT.

110 1.1 Anomalies in SCADA

111 Anomalies are defined as instances that do not fit the patterns of the rest of the data.
112 This misfitting makes it appear that the instances in question have been generated
113 by an altogether different mechanism to that generating the "normal" data [3]. It
114 is important that these anomalies are identified and removed so as to not bias the
115 relationships being studied. For WT *SCADA* data, anomalies are typically categorized
116 into 3 types [4] and any anomaly detection (AD) method must be designed with these

117 in mind. These are described below and illustrated with **Figure 1**:

- 118 • Anomaly type 1: These anomalies are characterized by no power output whilst
119 above cut in wind speed i.e. parking/downtime imposed by the operator. In
120 general AD terms, these instances would be characterized as *contextual* anoma-
121 lies [3]. An instance is a contextual anomaly if it would be considered normal
122 in a different context. The values of each feature are normal in isolation, but
123 abnormal when considered together.
- 124 • Anomaly type 2: These anomalies are characterized as steady and continuous
125 positive power output at a power less than the turbine’s rated power, P_r , i.e.
126 curtailment imposed by the operator. These can be characterized as *contextual*
127 anomalies as well.
- 128 • Anomaly type 3: These anomalies are randomly scattered across the feature
129 space. Reference [4] suggested these may be caused by sensor malfunction or
130 noise in signal processing. It is also possible these instances are generated by
131 stop-to-operation transitions or vice versa. These are best described as *point*
132 anomalies. Point anomalies are single instances that clearly do not conform to
133 the nature of the rest of the data.

134 Note that "Contextual" anomalies are exactly that, driven by context so will dif-
135 fer depending upon the end use of the SCADA data. For applications that seek to
136 understand turbine performance, we would consider any significant deviation from
137 the manufacturer’s specified WTPC as anomalous. For example, an instance where
138 no power was generated at wind speeds above cut-in, such as curtailment by the op-
139 erator, we would consider an anomaly. For condition monitoring applications, the
140 same curtailed instance would not be considered anomalous as the operating mode is
141 included in the context. This paper focuses on applications concerned with turbine
142 performance so considers deviation from the WTPC as anomalous, be it from operator
143 influence or otherwise.

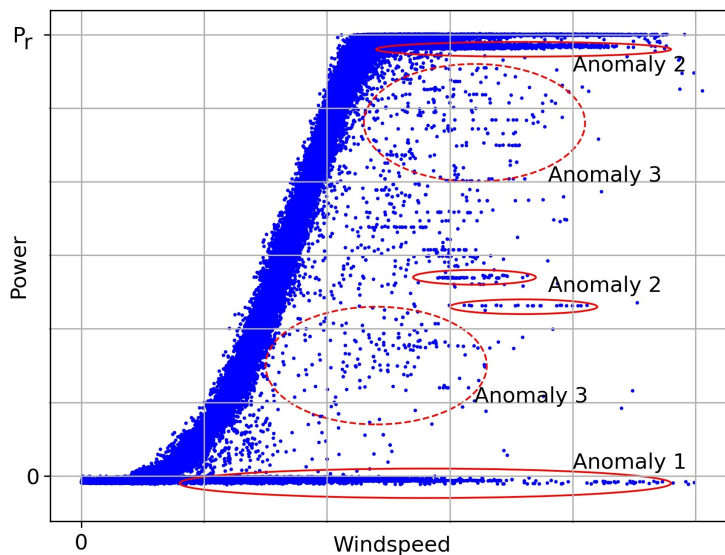


Figure 1: Wind turbine power curve showing examples of anomaly types. P_r indicates rated power of the WT.

144 To emphasize, contrary to the bulk of AD research, the mechanisms of anomaly
 145 types 1 and 2 are well understood. Fault, downtime, and curtailment instances should
 146 be explicitly labelled in the SCADA by any competent system. It only remains for the
 147 user to remove these instances.

148 1.2 Literature Review

149 Taking a broad view, the literature can be split into two groups based on their ap-
 150 proaches to cleaning WTPC data: group ①, those that pre-process data prior to
 151 running AD and group ②, those that do not. This is, of course, among many divi-
 152 sions that can be made.

153 In group ②, the non-pre-processing group, a wide range of approaches have been
 154 formulated. Notable examples include defining normal behaviour [5], using image-
 155 based approaches [6], using a "Change Point Grouping Algorithm and Quartile Algo-
 156 rithm" [7], and using Gaussian Mixture Models (GMM) [8]. Common themes from
 157 these works is the difficulty in dealing with "stacked" data, this is where anomalies

158 are so numerous they start to sway the statistical perception of "normal". Another
159 theme is the difficulty to reliably produce good results when anomalies form increas-
160 ingly large percentages of the data. All of these papers go straight to the anomaly
161 detection method with no pre-processing of data to remove obvious outliers or missing
162 data.

163 In contrast to the above, group ① papers, those that implemented pre-processing,
164 appear to have approached WTPC cleaning from a data-mining standpoint. In [9]
165 the authors pre-processed SCADA data into categories of "unnatural", "constant",
166 "exceeding", "missing" or otherwise valid data. This was followed by determination of
167 "irrational" data using a 2-step process. First, each instance of the remaining data was
168 given a weighting depending on its distance from the manufacturer-specified WTPC.
169 Second, the Local Outlier Factor (LOF) technique was applied with these weightings.
170 The entire process resulted in some 4,190 instances of a 18,001 dataset being removed
171 in their case study. Such explicit use of the manufacturer's specifications, i.e. cleaning
172 based on how it performs in theory rather than in practice, is completely at odds with
173 group ②'s unsupervised approach. The approach is questionable given that it is not
174 straightforward to compare the specified WTPC of a WT to that being achieved on a
175 given site. To point out the most glaring obstacles, differences in topography, terrain
176 roughness, and wind regime will need to be compensated for, as will any potential
177 wake effects from nearby turbines. There is a risk of introducing bias into the data,
178 a non-problem of the unsupervised group ②. Aside from this, [9] dealt with obvious
179 invalid SCADA data, such as missing data, of which none of group ② papers even
180 mention. Logically, any AD technique will find defining "normal" easier if valid data
181 makes up a greater percentage of the dataset; [9] pointed out that the LOF, which
182 uses distances between instances, would struggle with stacked data. Unfortunately,
183 the impact on power curve cleaning with and without pre-processing is not compared
184 by the authors.

185 Another paper that employed pre-processing was [10]. This study took a multi-step
186 approach and implemented simple statistical methods prior to running the DBSCAN
187 (density based spatial clustering of applications with noise) technique. The first step

188 eliminated negative power instances. The second and third steps applied the box plot
189 rule (see **Section 2.4** for an explanation of the box plot rule) to wind speed and
190 power intervals respectively. Finally, the DBSCAN method was applied. The authors
191 stated that the purpose of applying the box plot rule was to eliminate sparse outliers,
192 so making the boundaries of the stacked outliers clearer and improve the efficacy of
193 DBSCAN. The impacts of steps 1-3 were, unfortunately, not evaluated.

194 Of all the papers referenced in the literature review, only two papers dealt with
195 missing data, references [9] and [11]. The latter study, [11], concerns the detection of
196 blade icing and does so by comparing 3 methods. These are: percentage deviation from
197 the manufacturer’s WTPC, standard deviation of power for a wind speed interval, and
198 using quantiles of power for a wind speed interval. Furthermore, this paper is the only
199 one to use explicit fault or curtailment indications in the SCADA data; however, it
200 should be noted that [9] did reference a study that used operator logs. This general
201 lack of acknowledgment is unusual given that the international standard for power
202 curve measurement, IEC 61400-12, prescribes a ”data quality check” of removing
203 ”unavailable” or ”out of range” measurements and data rejection based on power
204 limited instances and faults with referencing to operator logs [12].

205 On the topic of ice detection, deviation from the WTPC is a common approach
206 for detection of ice accumulation on the blades. The authors note that ice detection
207 studies are rarely referenced by the general WTPC cleaning community, and vice
208 versa. Icing typically manifests itself as a deviation from the WTPC, more specifically
209 a reduction in power compared to manufacturer’s specification. Other studies on this
210 topic include [13] and [14]. In [13], the supervised learning Random Forest classifier
211 (not to be confused with iForest) is used on pre-processed and labelled data. The data
212 was pre-processed by the WT operator prior to being handed over to the authors,
213 therefore the precise pre-processing methods are not discussed. In [14], kNN-regression
214 is employed; however, data pre-processing is not mentioned.

215 The prevailing attitude of group ② appears to be the less supervision a technique
216 requires the better, provided this is not at the expense of results. Group ① better
217 embraced the spirit of SCADA data but rarely acknowledge that explicit or implicit

218 indications of faults or curtailment exist, or suggest actively using this knowledge
219 as part of pre-processing. Hesitation around using these indications is understandable
220 given the potential for mislabelling, lags between a fault occurring and the alarm being
221 logged, or the possibility that these logs simply do not exist. It follows that if an AD
222 method can perform just as well without these logs then it would be undesirable to
223 add steps into the process to unnecessarily act upon these logs. However, given the
224 lack of acknowledgement, it is hard to say whether this hesitation is justified. It is the
225 author's experience that these explicit indications of operating status are included in
226 SCADA far more regularly than they are not. After all, SCADA is Supervisory Control
227 and Data Acquisition, if a system does not record these parameters it can hardly be
228 called a SCADA system. As such, the lack of acknowledgement in the literature is
229 surprising. The choice to omit knowledge of indications from an AD method is still a
230 choice.

231 **1.3 Contributions and Paper Organisation**

232 The key contributions of this paper to current knowledge gaps are as follows:

- 233 • The impact of filtering SCADA of explicit and obvious anomalies, such as faults
234 and curtailment, based on indications in the data prior to running AD techniques
235 was found to be understudied in the literature. This paper investigates this
236 by comparing the performances of AD techniques with and without filtering
237 applied.
- 238 • Going beyond simply filtering out obvious/explicit anomalies, this paper also
239 explores utilizing this data further. Where possible, the filtered data is used
240 to form an anomaly-class. This allows classification-based anomaly detection
241 methods to be used. Such an approach has not been found in the literature.
- 242 • Underpinning the previous 2 contributions, simple rules for the filtering are
243 developed based upon explicit fault and curtailment indications in the SCADA,
244 as well as pitch values. These rules avoid the need to overly specify expectations
245 of how the turbine should be performing, i.e. avoids the need to incorporate the

246 manufacturer-specified WTPC.

247 • With the large assortment of AD techniques present in the literature, this paper
248 provides comparison between 4 distinct techniques on their ability to clean the
249 WTPC, whilst still maintaining the statistical variability in the wind speed
250 feature.

251 • Whilst papers in the literature have previously evaluated AD techniques on
252 their ability to maintain data, ability to maintain statistical variability in the
253 wind speed feature appears understudied. Such a metric would highlight good
254 performers. This metric is introduced and applied here.

255 • Finally, the importance of proper treatment of missing data, that is to say
256 instances with some or all values missing or corrupted, is raised. Treatment of
257 missing data appears to be rarely mentioned in the literature, in this paper the
258 importance of such data in the context of anomaly detection is discussed.

259 The remainder of this paper is organized as follows. Section 2 describes the
260 methodology used in this paper. This includes descriptions of the 4 AD techniques,
261 their operation, and how they are applied. Additionally, descriptions of 3 different
262 filtering-based approaches to SCADA data are given. Section 3 describes the WTs
263 and wind farms the SCADA data used in this paper originate from. Section 4 details
264 and interprets the impacts of each AD method as well as comparing between them.
265 Finally, Section 5 summarizes the key findings and contributions of the paper.

266 2 Methodology

267 The methodology of this study is shown in **Figure 2**. The methodology is composed
268 of 3 main components:

- 269 • Missing data treatment.
- 270 • Anomaly detection and removal
 - 271 – Explicit and obvious anomaly filtering.
 - 272 – Anomaly detection proper.

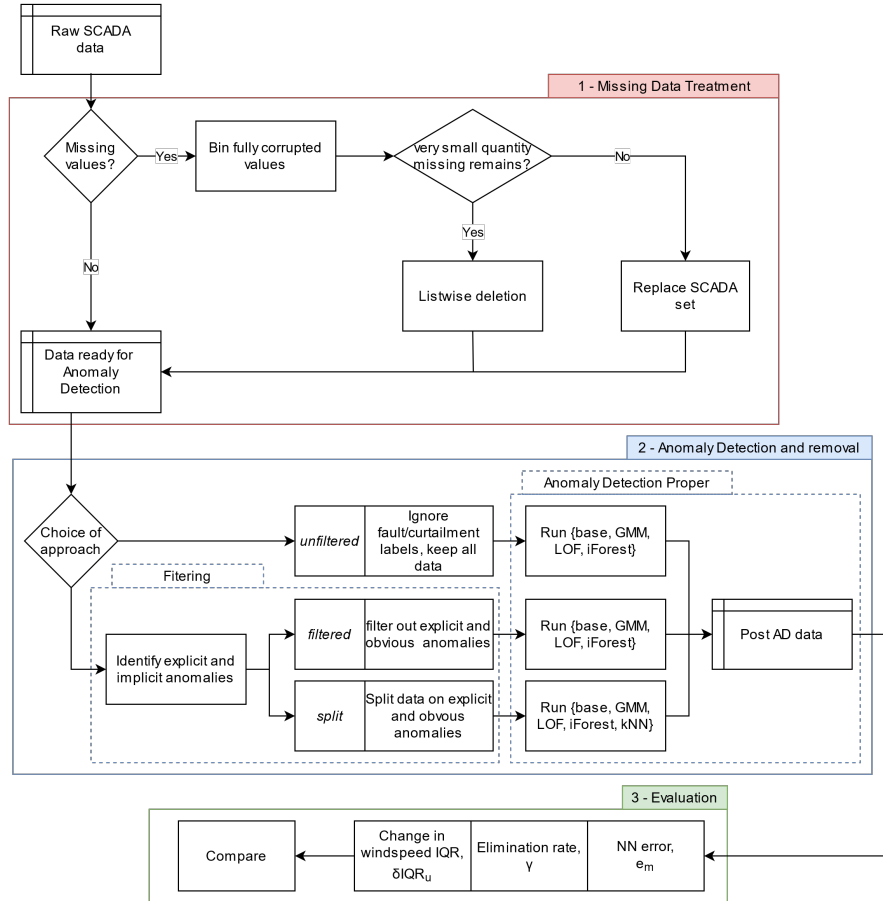


Figure 2: Methodology for treatment of SCADA data for anomaly detection.

274 Treatment of missing data is discussed in **Sections 2.1**. From the choice to filter or
 275 not, 3 approaches of *Unfiltered*, *Filtered*, and *Split* are proposed in **Section 2.2** along
 276 with a description of how filtering is performed. Five anomaly detection techniques,
 277 including "do-nothing", are described in **Section 2.3** along with how they are applied.
 278 The combination of an approach and an AD technique is referred to herein as an "AD
 279 method". Finally, the evaluation of each AD method is detailed in **Section 2.4**.

280 2.1 Missing Data Treatment

281 The treatment of missing data appears to be rarely discussed in papers relying on
282 SCADA data. Looking again at the papers, only [9] and [11] make reference to handling
283 missing data. As discussed in [15], missing data can be categorized into 3 groups:

- 284 • Missing Completely at Random (MCAR): There is no correlation between the
285 missingness of data and any variables.
- 286 • Missing at Random (MAR): Missingness is correlated with a variable. The cause
287 can be measured and included in missing data methods. Failure to include this
288 would introduce bias into future data models.
- 289 • Missing Not at Random (MNAR): Like MAR, missingness is correlated with a
290 variable. However, the cause cannot be measured and so cannot be corrected
291 for in missing data methods.

292 The author notes the parallels between *missing* data and *anomalous* data. Take,
293 for example, MAR and anomaly type 3. The descriptions are almost interchangeable.
294 Reference [15] notes that common approaches to missing data are listwise deletion and
295 data imputation. Listwise deletion is when an instance is simply deleted entirely. Both
296 methods have been considered outdated for some time now due to their potential to
297 introduce bias into data models. As a precursor to AD, one can see how these methods
298 might be counterproductive by changing the nature of the data models.

299 Keeping in mind the aim of not introducing bias, listwise deletion is only appro-
300 priate if the data to be removed represents a small percentage of the overall data.
301 This not only avoids introducing bias, but also maintains statistical power. Reference
302 [15] suggested that less than 5% would be trivial. If greater percentages exist, then
303 one must determine if the missingness is MCAR, MAR, or MNAR. For the purposes
304 of this document, this is typically tested by a software package. Should MCAR or
305 MAR be the cause, then *Full-Information Maximum Likelihood* (FIML) procedures
306 are recommended. Similarly, the data imputation that would follow is also handled
307 by the software and the theory is not explored here. In the unlikely event that the
308 type of missingness is MNAR, then the methods available to correct the data become

309 few. Reference [15] suggested use of *Pattern-Mixture methods*, however, from a prag-
310 matic approach to SCADA data WTPC cleaning, further use of that dataset should
311 be questioned.

312 In the context of WT SCADA data, it is common to have many instances that are
313 entirely empty, barring features for WT identification, etc. These can be generated
314 by many means such as clock time changes. These instances do not fall within the 3
315 missing data categories and should be removed entirely.

316 2.1.1 Application of Missing Data Treatment

317 As per **Figure 2**, fully corrupted instances are first removed. The number of instances
318 containing missing data is then calculated. If these represent less than 5% of the data
319 then they are considered trivial (as per [15]) and are removed via listwise deletion.
320 Given the number of WT SCADA sets available, if this 5% threshold is exceeded, the
321 SCADA set will be abandoned and another adopted. This is a pragmatic approach to
322 avoid the need to employ FIML software.

323 2.2 Explicit and Obvious Anomaly Filtering

324 As discussed in **Section 1.2**, SCADA data, by its nature, contains explicit indications
325 of the operational state of a WT, such as fault alarms and operational time. Fur-
326 thermore, some instances can be characterised as "obvious" anomalies, namely high
327 blade pitch values indicate the operator was stopping the turbine. Similarly, power
328 reference values state the maximum power limit the operator imposed on the WT for
329 each period. A value less than P_r indicates curtailment. These anomalies can be fil-
330 tered from the datasets prior to running AD techniques. Alternatively, all data could
331 be maintained, or the filtered data could be further leveraged. Arising from this, 3
332 distinct approaches are proposed:

- 333 • *Unfiltered*: No filtering occurs. As such, all SCADA data is kept and fed into
334 AD processes.
- 335 • *Filtered*: Filtering occurs and the filtered data is removed entirely. The remain-

336 ing data is fed into AD processes. Issues of stacking and high initial anomaly
337 percentages should be reduced. However, any mislabelled instances are lost from
338 the dataset, along with statistical power.

- 339 • *Split*: The dataset is split in two according to the filtering rules. This forms an
340 assumed-anomaly set and an assumed-normal set. AD techniques are run using
341 both sets to identify instances that are mislabelled (i.e. in the wrong set). This
342 approach attempts to make most use of all data.

343 Comparison of *Unfiltered*, *Filtered*, and *Split* approaches will identify the impact
344 of pre-processing the SCADA data.

345 2.2.1 Application of Filtering

346 SCADA data will be filtered out for *Filtered* and *Split* approaches if any of the following
347 conditions are met:

- 348 • SCADA "Fault" feature indicates a fault;
- 349 • SCADA "Operational time" feature indicates the WT was not in operation for
350 the full duration of recording period. In practice, this is if the WT operated for
351 less than 10-minutes in the 10-minute period.
- 352 • SCADA "Power reference" feature value is less than 99% of P_r . This indicates
353 the WT was curtailed.
- 354 • SCADA blade pitch angle feature value is greater than 30° . This indicates the
355 WT was stopped by the operator. As shown in **Figure 3**, stoppage occurs at
356 pitch angles of $80-90^\circ$, a threshold of 30° is chosen to capture some operation-
357 to-stop transition instances whilst not incorrectly filtering out normal instances.
358 From **Figure 4**, it is clear that normal operation pitch angles end at approxi-
359 mately 25° .

360 RPM will not be used in indicating curtailment. Whilst a lack of rotation above
361 cut-in wind speed would likely indicate curtailment, this introduces the need to define
362 a WTPC. Given that the mechanism for stoppage is blade pitching (and braking),

363 using pitch avoids the need to define the WTPC or normal behaviour. This saves a
364 considerable amount of time and effort for the user.

365 Typical results of the filtering process are shown in **Figure 3**.

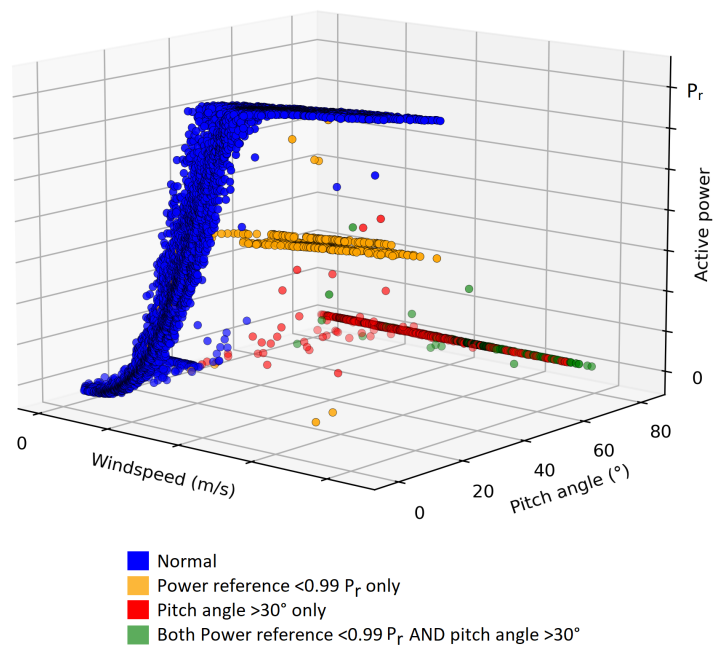


Figure 3: Visualisation of the impacts of explicit and obvious filtering applied to SCADA data from a WT from wind farm B. As shown by the orange, red, and green instances, large amounts of non-normal instances can be filtered out. The blue instances scattered across the figure, as well as the cluster of instances around 25° pitch at 0kW , show that filtering alone is not sufficient and further AD techniques need to be applied.

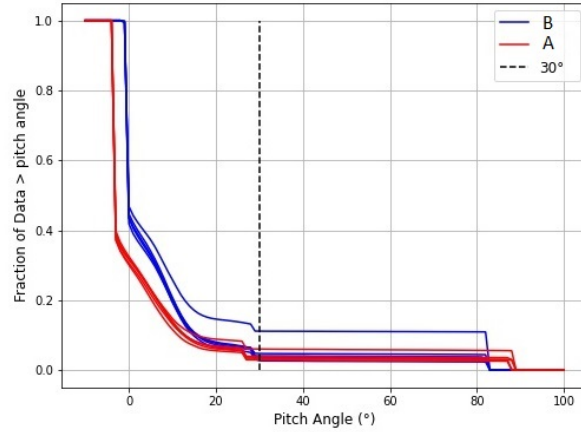


Figure 4: Plot of pitch angle versus fraction of dataset smaller than said pitch angle. Ten WTs from two wind farms are shown. A pitch angle of 30° is indicated, greater than this an instance is labelled as curtailment according to the filtering rules.

366 2.3 Anomaly Detection Algorithms

367 Five AD techniques have been chosen, including "do-nothing". Each method is named
 368 using the format *approach.technique* and all method names are summarized in **Table**
 369 **1**. Their theory and how they are applied is described below.

Table 1: Names of AD methods used in analysis.

Technique	Approach		
	Unfiltered	Filtered	Split
base (none)	unfiltered.base	filtered.base	split.base
iForest	unfiltered.iForest	filtered.iForest	split.iForest
GMM	unfiltered.GMM	filtered.GMM	split.GMM
LOF	unfiltered.LOF	filtered.LOF	split.LOF
kNN	N.A.	N.A.	split.kNN

370 Where data is scaled, this will be performed using robust scaling. A robust scaler
 371 is used under the assumption that the data contains outliers. Scaling is performed via
 372 *Scikit* module *preprocessing.RobustScaler*[16]. This is a standardization method which
 373 scales data using the interquartile range (IQR, see **Section 2.3.3** for a description of
 374 IQR). As such, it is more robust to outliers than Min-Max scalers. Min-Max scalers

375 scale all values between 0 and 1. A spuriously large values value, say a reading erro-
376 neously 10 times larger than the true value, would lead to all the other values being
377 crushed into a small range during Min-Max scaling.

378 **2.3.1 Do Nothing - Base**

379 This "technique" represents no action being taken. Note that this is still after filtering,
380 hence the *unfiltered.base* method uses all data and the *filtered.base* method uses only
381 filtered data. The *split.base* method is identical to the *filtered.base* method as, by the
382 nature of this being the base method, no further action can be taken.

383 **2.3.2 iForest**

384 Isolation Forest (iForest) is a relatively new technique in the field of AD and was
385 developed and introduced by [17]. iForests are ensembles of "iTrees". These iTrees are
386 an evolution of binary search trees in that they partition data, however, in iForests
387 the splits are made at random. For multiple features, a random feature is selected
388 followed by the random split. The shorter an instance's path length the more likely
389 an instance is to be an anomaly. The underlying theory being that normal instances
390 occur in the same region as other instances, hence requiring many splits to be isolated.
391 Conversely, anomalies exist in sparsely populated regions and so require far fewer splits
392 to be isolated. This is shown in **Figure 5**.

393 iTrees work through the dataset in samples, rather than process the entire dataset
394 in one step. According to the authors, a small sample sizes allows iForest to overcome
395 problems of masking and swamping. A sample size of 2^8 (256) is recommended.

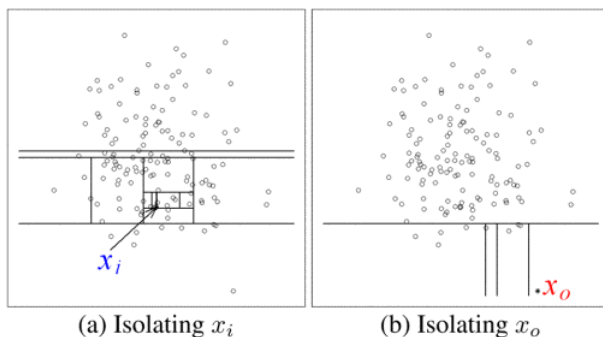


Figure 5: Visualisation of the iForest technique isolating a normal instance (a) and an anomalous instance (b). Figure from [17].

396 As per the original paper, iForests derive the anomaly score for an instance x from
 397 its averaged path length, $h(x)$. The anomaly score for an instance x , given a set of n
 398 instances, is given as:

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad (1)$$

399 Where $c(n)$ is the average of $h(x)$ given n and $E(h(x))$ is the average path length of
 400 x across all the iTrees. If s is close to 1, the instance is likely to be an anomaly, if less
 401 than 0.5, then it is likely not an anomaly. According to the authors, if all instances
 402 return an s value of approximately 0.5 then the set does not have any anomalies.

403 The main advantage of iForests are that they are extremely quick to run. iForest
 404 does not perform any profiling, distance, density, or co-variance calculations so the
 405 computational power required is tiny relative to other AD techniques. iForest process-
 406 ing time can be further reduced by imposing a height limit upon the iTrees, beyond
 407 this an instance would be considered normal.

408 iForest will be implemented in Python using *Sklearn-ensemble-IsolationForest*, an
 409 algorithm from *Scikit-Learn* [16]. Default settings of 100 trees and a sample size of 2^8
 410 will be used, along with no assumption for percentage of contamination. Note that no
 411 data scaling will be required prior to applying the iForest technique. With reference to
 412 the *Split* approach, there appears to be no way to further utilize the initially-anomalous
 413 data. As such, *split.iForest* and *filtered.iForest* methods are identical.

414 2.3.3 Gaussian Mixture Modelling

415 Gaussian Mixture Model (GMM) is a model-based technique that can be adapted
416 for AD. The underlying assumption is that the model being analysed is composed
417 of k Gaussian distributions. Normal instances are generated from these Gaussian
418 distributions whilst anomalies are not and so occur in low probability spaces [3].

419 As per [18], for p variables, each distribution has a mean for each variable, $\mu =$
420 $(\mu_1, \mu_2, \dots, \mu_p)$. Each distribution will also have a covariance matrix, containing covari-
421 ance values for each pair of variables. The means and covariance matrix values are
422 estimated using Maximum Likelihood Estimates.

423 Anomaly scores are simply the distance from the instance to the mean. This
424 is usually Euclidean distance, however some methods use Mahalanobis distance. As
425 such, each instance has as many anomaly scores as there are Gaussian distributions
426 assumed. There appears to be no definitive way to convert these anomaly scores to
427 classifications but many have been proposed. Reference [19] suggests assigning any
428 score greater than 3 standard deviations away from the mean score as an anomaly.
429 Reference [20] suggests using the *box-plot rule*. This is the range between the whiskers
430 of a box plot, equivalent to p_{lower} and p_{upper} given as:

$$[p_{lower}, p_{upper}] = [Q_{25} - 1.5 \times IQR, Q_{75} + 1.5 \times IQR] \quad (2)$$

431 Where IQR is the interquartile range, equivalent to the difference of Q_{75} and Q_{25} .

432 The number of mixtures used in the model is typically determined using BIC
433 (Bayesian Information Criterion) curves. The theories behind BIC curves are not
434 explored here. The number of mixtures producing the lowest BIC scores should be
435 used, however, this is not an absolute rule and often the curve "elbow" will be used.
436 The elbow is, essentially, the point of diminishing returns. Adding more mixtures no
437 longer results in a similar drop in BIC score as it did for adding mixtures previously,
438 even if some small drop is witnessed [21].

439 GMM will be implemented in Python using `Sklearn.mixture.GaussianMixture`, an
440 algorithm from `Scikit-Learn` [16]. The number of mixtures to use has been determined

441 using the BIC curve process on a reduced set of features. These were: all wind speed
 442 features, temperature, and mean active power (see **Section 3** for full description of
 443 SCADA features). Ten randomly selected turbines were selected (5 from each wind
 444 farm, see **Section 3**) and their curves determined. The curve elbow was found to be
 445 at 3 mixtures as shown in **Figure 6**.

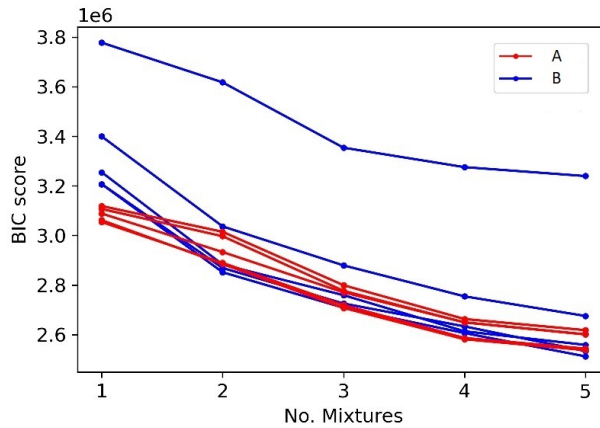


Figure 6: BIC scores versus number of Gaussian mixtures assumed for 10 randomly selected turbines (5 from each wind farm). Raw SCADA data was used.

446 Unlike the other techniques, *Sklearn-mixture-GaussianMixture* does not have a
 447 build-in classification method; however, it does have the *score_samples* method, which
 448 calculates probabilities of each instance belonging to each of the 3 mixtures. For each
 449 instance, the maximum of the 3 likelihoods will be taken. The box-plot rule will then
 450 be applied to determine normal and anomalous data.

451 In the case of *unfiltered.GMM* all data will be used. For *filtered.GMM* only filtered
 452 data will be used. For *split.GMM*, an altogether different approach will be taken. The
 453 SCADA sets will be split into 2 groups based on the filtering, an assumed-normal
 454 group and an assumed-anomalous group. From each group, the data will be further
 455 split into train and test data on an 80-20% split. The training data will be used to
 456 train a GMM per group. The testing data will be fed into **both** models to see which
 457 model they have greater affinity. The results are recorded. This entire process is

458 repeated 4 more times with different portions of the data forming the test data (i.e.
 459 cross-validated) until every instance has been assigned as being more likely normal or
 460 more likely anomalous.

461 2.3.4 LOF

462 Local Outlier Factor (LOF) is a relative-density based technique proposed by [22].
 463 LOF is concerned with assigning a degree of outlier-ness to an instance, rather than
 464 a hard label of 0 or 1. The LOF method uses a local approach, rather than global,
 465 hence it is useful in applications where non-homogeneous densities are acceptable.

466 The LOF method first requires the number of nearest neighbours, $MinPts$, to be
 467 chosen. Following this, the LOF of a point x is given as:

$$LOF_{MinPts}(x) = \frac{\sum_{o \in N_{MinPts}(x)} \frac{lrd_{MinPts}(o)}{lrd_{MinPts}(x)}}{|N_{MinPts}(x)|} \quad (3)$$

468 To interpret the above, the LOF of an instance is the ratio of that instance’s local
 469 reach density, lrd , to that of the average lrd of its $MinPts$ neighbours. The lrd of
 470 instance x is given as:

$$lrd_{MinPts}(x) = 1 / \left[\frac{\sum_{o \in N_{MinPts}(x)} reach-dist_{MinPts}(x, o)}{|N_{MinPts}(x)|} \right] \quad (4)$$

471 The $reach-dist$ is essentially the euclidean distance between two points, however,
 472 if the distance becomes very small, this reverts to a set value. This minor change
 473 has the effect of smoothing LOF scores and making differentiation between inliers and
 474 outliers easier.

475 In **Figure 7** the lrd of an anomalous instance p with $MinPts$ value of 3 is illus-
 476 trated. The combined reach distances of p to its 3 nearest neighbours is clearly far
 477 larger than the same metric for those 3 neighbours. This results in a low lrd compared
 478 to the neighbours, i.e. a large volume is required to capture the specified nearest
 479 neighbours. Such a large relative difference in lrd then results in a high value for
 480 $LOF_{MinPts}(p)$, as per **Equation 3**.

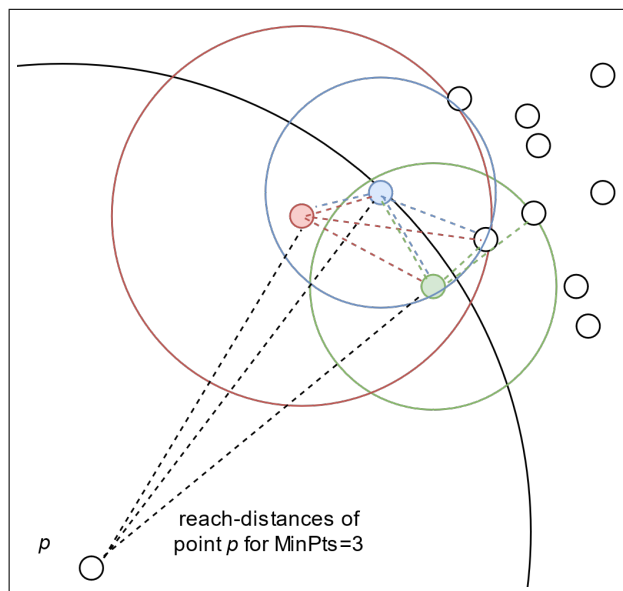


Figure 7: LOF in action, adapted from [22]. The 3 nearest neighbours of p are shown by the green, blue, and red points. Each of their 3 nearest neighbours are shown with the dashed lines of the same colours. Instance p is clearly an outlier with respect to the rest of the set due to its large reach and resulting low local reachability density.

481 According to [22], most instances should have a LOF value of 1. If the LOF value
 482 is greater than 1, it is likely an outlier and sits in a sparsely populated region. A LOF
 483 of less than 1 and it is likely an inlier in a densely populated region. There appears
 484 to be no definitive threshold LOF value to label an instance an outlier. The value
 485 that will be used in this study is 1.5, the default used by the Python module *sklearn-*
 486 *neighbors-LocalOutlierFactor*, as discussed below. Similarly, there is no single value for
 487 *MinPts* recommended by the authors of the LOF method, although no smaller than
 488 10 is recommended to avoid statistical fluctuations. An appropriate value is dependent
 489 on the dataset being investigated, if too large a value is used then clusters that are
 490 small, but valid, are unfairly treated.

491 LOF will be implemented in Python using *sklearn-neighbors-LocalOutlierFactor*, an
 492 algorithm from *Scikit-Learn* [16]. A *MinPts* value of 700 will be used. This is based
 493 upon a brief analysis of *MinPts* with removal rates for randomly selected turbines from

494 both wind farms. Removal rates were found to settle (in terms of visual impact on
495 the WTPC) at approximately $MinPts = 400$, with no change between this value and
496 $MinPts = 1000$. It may seem like an unfair advantage has been gotten by choosing
497 the most favourable value of $MinPts$; however, note that all the data is unlabelled
498 and that choice of $MinPts$ here is based on removal rates, not accuracy. The fact that
499 a conservative value of $MinPts$ was "manually" found is a matter of convenience to
500 save coding and computational time, as this process could be automated.

501 It is important to note that multidimensional data must be scaled prior to use in
502 LOF so as to account for the different dimensions of features in distance calculations.

503 LOF has no training stage and has only a fit-predict method. Note that the
504 nature of LOF is to make an astronomical number of calculations. As such, for high
505 dimensional data, users of this technique may need to reduce the number of dimensions
506 used to make the computation time feasible. With reference to the *Split* approach,
507 there appears to be no way to further utilize the initially-anomalous data. As such,
508 *split.LOF* and *filtered.LOF* methods are identical.

509 **2.3.5 k Nearest Neighbours**

510 k Nearest Neighbours (kNN) is a nearest-neighbours technique. This technique is a
511 classification method but is adapted here for AD. kNN dates back to 1951 and was
512 introduced in [23]. The underlying assumption is that normal instances occur close
513 to other normal instances, and the same is true for anomalies. As such, this method
514 requires that labelled instances already exist so that new, unlabelled, instances can be
515 classified.

516 The kNN technique first involves specifying the number of nearest neighbours to
517 consider, k . A distance metric, such as Euclidean, is then specified. An unlabelled
518 instance is then considered against a labelled dataset. Its k nearest neighbours are then
519 determined. The unlabelled instance is then assigned to whichever class comprises the
520 majority of neighbours. This is illustrated in **Figure 8**. Variations upon this technique
521 can be to add weightings to certain instances, for example if these instances have been
522 chosen by experts.

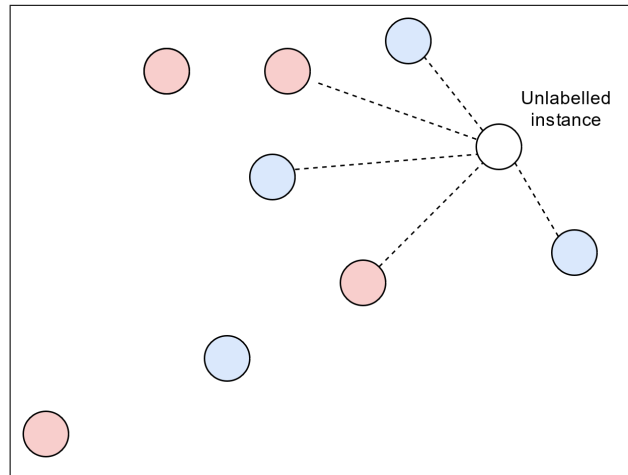


Figure 8: Assigning an unlabelled instance a class via kNN using $k = 5$. The new instance would be assigned to the blue class.

523 kNN will be implemented in Python using `Sklearn.neighbors.KNeighborsClassifier`,
 524 an algorithm from *Scikit-Learn* [16]. This technique is exclusive to the *Split* approach
 525 as it requires classes to be assigned. First, the best k value will be chosen from options
 526 of 3, 5, 7, and 9. This will be determined by assigning all data classes based upon
 527 filtering. Then, the data will be split into train and test data on an 80-20% split. The
 528 value of k that produces the highest accuracy for the test set will then be used.

529 *Split.kNN* will then be implemented in a similar fashion as *split.GMM*. All data
 530 will be given labels based upon filtering. Eighty percent of each group (normal and
 531 anomalous) will then be used to train a classifier. The classes of remaining remaining
 532 20% will then be predicted using the k value determined above. This will be repeated
 533 4 further times, with different portions of the data forming training and testing sets,
 534 to determine a classification for every instance, i.e. 5-fold cross validation.

535 Note that there is a need to scale features when using kNN due to the use of
 536 distance metrics.

537 **2.4 Evaluation**

538 To paraphrase from [8], the goal of AD in cleaning the WTPC is "to reject potential
 539 [anomalies] whilst broadly retaining the statistical characteristics of the WTPC, in
 540 particular the mean values of the measurements". As such, the efficacy of each AD
 541 method is assessed by three measures:

- 542 • Prediction error on the "cleaned" dataset, e_M .
- 543 • Elimination rate, i.e. percentage of data removed, γ .
- 544 • Change to wind speed interquartile range, δIQR_u .

545 **Prediction error**, e_M : A neural network (NN) will be fitted to each AD method's
 546 cleaned WTPC using a standard train-test split. The prediction error will then be
 547 found and converted to a percentage of the P_r to allow for comparison between turbines
 548 of different ratings. This is calculated as follows:

$$e_M = \frac{RMSE(u, h)}{P_r} \times 100\% \quad (5)$$

549 Where P_r is the rating of the WT, and $RMSE(u, h)$ is the test-set prediction root
 550 mean square error of a neural network trained on the training set, given by:

$$RMSE(u, h) = \sqrt{\frac{\sum_{i=1}^n (h(u_i) - p_i)^2}{n}} \quad (6)$$

551 Where n is the total number of instances in the prediction set, $h(u_i)$ is the predicted
 552 value of power for instance i with wind speed u , and p_i is the actual power for that
 553 instance.

554 A simple NN will be constructed using *TensorFlow* and run using *Google Colab's*
 555 TPUs (Intel® Xeon® CPU @ 2.30GHz, 64 GB RAM). A shallow configuration is
 556 chosen to allow for testing within a reasonable time frame. With 10 unique methods
 557 to run across 20 WT SCADA sets, this results in constructing, training, and testing
 558 200 NNs. The NN will consist of 1 input layer, 1 output layer, and 2 dense inner layers
 559 each with 64 neurons. Activation functions of the middle layers will be Rectified Linear
 560 Units. Given that the NN has only 1 input (wind speed) and 1 output (power), and

561 that the focus of results is in the comparison between methods, such a shallow NN is
 562 deemed appropriate and optimization of each NN is unnecessary. Obtaining results for
 563 a large number of turbines, therefore determining more reliable averages, is preferable
 564 to highly accurate results for a small numbers of turbines.

565 **Elimination rate**, γ : As per [10], elimination rate is defined as the percentage of
 566 data removed:

$$\gamma = \frac{N_b - N_a}{N_b} \times 100\% \quad (7)$$

567 Where N_b is the original number of instances in the SCADA set, and N_a is the
 568 number of instances after the AD method has removed the anomalies it has identified.
 569 If methods have similar prediction errors, the method that retains more data would
 570 clearly be superior. Conversely, any methods that rig the system by removing high
 571 percentages of data, so as to only make predictions over a very small range, can be
 572 identified as performing poorly.

573 **Change to wind speed interquartile range**, δIQR_u : The spread of the wind
 574 speed feature will be recorded both before and after anomaly removal via IQR. As per
 575 **Section 2.3.3**, IQR is simply the difference of Q_{75} and Q_{25} . Note: this is different
 576 from the box plot rule which includes the whiskers of the box plot. Change to IQR_u
 577 is calculated as:

$$\delta IQR_u = \frac{IQR_{u,b} - IQR_{u,a}}{IQR_{u,b}} \times 100\% \quad (8)$$

578 Where $IQR_{u,b}$ is the IQR of wind speed before the AD process, and $IQR_{u,a}$ is the
 579 IQR of wind speed after the AD process. IQR is chosen as it is robust to outliers than
 580 simply taking the absolute range as an anomalous measurement might read an order of
 581 magnitude higher, thus the change to absolute range would be drastic if this instance
 582 was (correctly) removed. By incorporating this IQR change and percentage removal
 583 a better representation of the method's quality can be achieved. Given that the AD
 584 methods will be applied in an unsupervised way, with error being calculated using a
 585 dataset the AD method has been applied to, in theory, a method could achieve a low
 586 error by reducing the SCADA set down to a narrow, predictable band. Incorporating

587 data and wind speed variability retention allows for better understanding of a method’s
588 appropriateness for AD in WTPC cleaning.

589 **3 Data Description**

590 **3.1 Wind Farms**

591 SCADA data from 2 offshore wind farms has been provided. These are referred to as
592 wind farms *A* and *B*. The characteristics of both wind farms are summarised in **Table**
593 **2**.

Table 2: Characteristics of the 2 wind farms used in the analysis.

Wind farm	A	B
Location	North Sea, UK	Northern Europe
Distance from shore	< 25km	< 25km
No. turbines	< 50	> 50
Turbines used in analysis	10, randomly selected	10, randomly selected
SCADA duration	24 months	18 months
SCADA frequency	ten-minute	ten-minute

594 **3.2 Feature Engineering**

595 The features provided in the SCADA sets is shown in **Table 3**. Note that wind speeds
596 are from nacelle anemometers for all WTs. Timestamps were not included in the AD
597 process, each instance is considered in isolation. All features of "Yaw" and "RPM"
598 were also dropped. These features are ineffective for determining curtailment (see
599 **Section 2.2.1**) and are not useful in generalizing the turbine power output for a site’s
600 given wind regime. Air density was not calculated due to a lack of pressure readings;
601 regardless, variations in air density have little impact upon the WTPC [24].

602 All other features are included in the AD process (including the treatment of
603 missing data). No further features were present in the SCADA set. The author notes
604 that SCADA features will vary between SCADA systems and some may contain more
605 explicit indications of operating mode.

Table 3: SCADA features with indication of use in Anomaly Detection. Blanks indicate that the feature was not present in the SCADA data.

Feature	Mean	Min	Max	Std. Dev.
Timestamp	N			
wind speed	Y	Y	Y	Y
Yaw	N	N	N	N
Pitch	Y	N	N	N
RPM	N	N	N	N
Power ref.	Y			
Power	Y	N	N	N
Temperature	Y			
Operation time	Y			
Fault	Y			

606 4 Results and Discussion

607 All SCADA datasets were cleaned of missing data instances via listwise deletion prior
 608 to performing filtering and AD. The percentage of the SCADA that was erroneous
 609 was less than 5% in every case and, therefore, less than the threshold value for being
 610 considered negligible. For wind farms *A* and *B*, erroneous instances made up an average
 611 of 0.23% and 0.06% respectively. This excludes instances in which all features were
 612 missing. As such, no substitution of SCADA sets was necessary.

613 Each anomaly detection method was applied to each of the 20 turbine SCADA sets.
 614 The impacts of each AD method upon the WTPC are visualized in **Figure 9** with
 615 colours indicating instances as being determined normal, anomalous, or filtered out
 616 by the AD method. Note that kNN technique only applies to the *split* approach and
 617 that iForest and LOF techniques are identical between *filtered* and *split* approaches,
 618 hence these results are not shown.

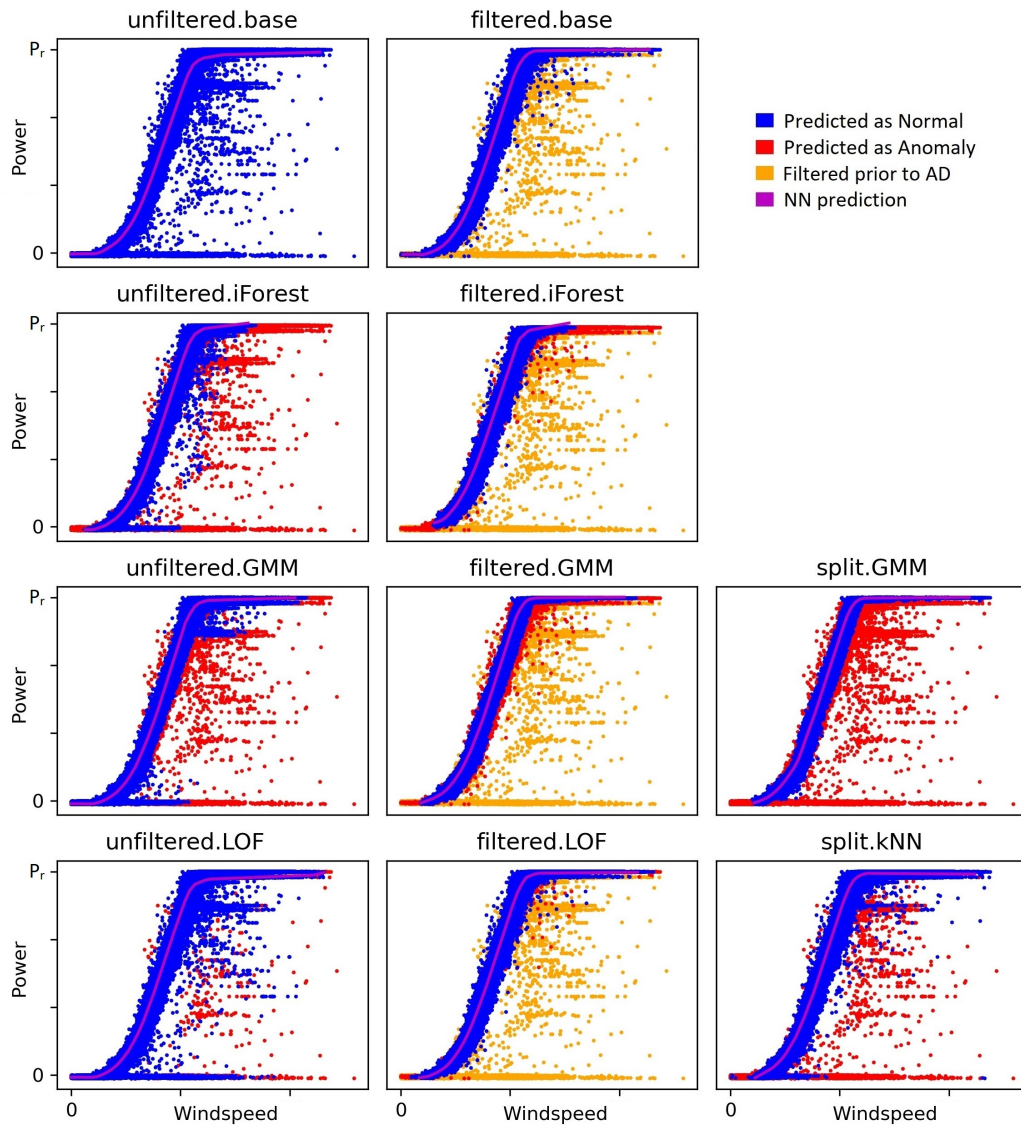


Figure 9: Visualisation of the impacts of different AD techniques and approaches when applied to a WT from wind farm *A*. As per the legend, different colours indicate whether instances were filtered out prior to AD or later labelled as "normal" or "anomalous" by the AD technique. In each subplot, the test-set predictions of the NN used for evaluation is shown by the purple line.

619 From an inspection of **Figure 9**, the differences in performances between tech-
 620 niques and approaches can be seen. Comparing *unfiltered.base* and *filtered.base*, it is

621 clear that applying simple filtering rules greatly cleaned the WTPC. However, many
622 Type 3 (random) anomalies still exist, as does a group of Type 2 (curtailment) in-
623 stances immediately under the flat part of the WTPC. This group of instances is most
624 visible in the LOF methods. Inspecting the *unfiltered* methods, it appears none man-
625 aged to deal with the Type 1 (stoppage) anomalies well, especially not at lower wind
626 speeds. Additionally, none managed to remove the Type 2 group discussed previously.

627 Looking at the iForest results, this technique did not handle the flat, P_r part of
628 the WTPC well. This technique does appear to have dealt with Type 3 anomalies
629 reasonably well. Moving on to the GMM technique, all three variations appear to
630 have removed the majority of Type 3 anomalies. *Unfiltered.GMM* appears to have
631 struggled at the knee of the WTPC and not removed a group of Type 2 anomalies
632 at approximately 80% of P_r . The *filtered* and *split* variations of GMM are the only
633 methods which removed the Type 2 group at approximately immediately under the
634 flat part of the WTPC. Inspecting the LOF technique, it appears to have been very
635 conservative and only removed the most isolated of instances. Stacking was clearly
636 an issue for *unfiltered.LOF*. *Filtered.LOF* only eliminated a handful of instances after
637 filtering occurred. Finally, kNN appears to have performed worse than *filtered.base*.
638 It appears initially anomalous points have been reclassified to normal, rather than the
639 other way around.

640 As per the methodology, prediction error as a percentage of P_r (e_M), elimination
641 rate (γ), and interquartile range of the wind speed feature (IQR_u) were determined
642 for each cleaned dataset, along with a percentage change (δ) from the base case, where
643 appropriate. The results have been averaged across the 20 turbines and are shown in
644 **Table 4** and **Table 5**.

Table 4: Results by method. Note that *split* and *filtered* approaches are identical for the techniques of base, iForest, and LOF. Elimination rate includes both filtered data and that detected as anomalous.

Method	e_M			γ	IQR_u		
	pre	post	δe_M (%)		pre	post	δIQR_u (%)
unfiltered.base	11.35	NA	NA	NA	5.86	NA	NA
unfiltered.iForest	11.35	4.89	56.95	16.02	5.86	5.16	11.94
unfiltered.GMM	11.35	5.53	51.23	6.20	5.86	5.77	1.60
unfiltered.LOF	11.35	8.06	28.94	1.67	5.86	5.84	0.47
filtered.base	11.35	3.78	66.71	10.97	5.86	5.54	5.50
filtered.iForest	11.35	3.73	67.17	27.47	5.86	4.78	18.51
filtered.GMM	11.35	3.38	70.25	15.13	5.86	5.50	6.29
filtered.LOF	11.35	3.72	67.19	11.74	5.86	5.52	5.91
split.GMM	11.35	3.39	70.11	14.11	5.86	5.48	6.62
split.kNN	11.35	4.05	64.35	7.99	5.86	5.62	4.12

Table 5: Subsequent rates of evaluation metrics, calculated by method. Note that *split* and *filtered* approaches are identical for the techniques of base, iForest, and LOF.

Method	$\delta e_M / \gamma$	$\delta e_M / \delta IQR_u$	$\delta IQR_u / \gamma$
unfiltered.base	NA	NA	NA
unfiltered.iForest	3.55	4.77	0.75
unfiltered.GMM	8.26	31.97	0.26
unfiltered.LOF	17.33	61.73	0.28
filtered.base	6.08	12.12	0.50
filtered.iForest	2.45	3.63	0.67
filtered.GMM	4.64	11.16	0.42
filtered.LOF	5.72	11.36	0.50
split.GMM	4.97	10.58	0.47
split.kNN	8.06	15.61	0.52

645 4.1 Impact of Filtering

646 The impact of filtering alone can be isolated by comparing *unfiltered.base* and *fil-*
647 *tered.base*. From **Figure 9**, it is clear that a large amount of data was removed, this
648 averaged 10.97% for these wind farms. This resulted in a substantial improvement
649 to prediction error, e_M , with an average decrease of 66.7% from *unfiltered.base* to *fil-*

650 *tered.base*. This occurred at a minimal change to IQR_u with only a 5% reduction.
 651 The average e_M of *unfiltered* and *filtered* methods was 7.46% and 3.65% respectively.
 652 Clearly, filtering was extremely beneficial to prediction error due to the reduced influ-
 653 ence of stacked anomalies.

654 4.2 Error and Elimination Rate

655 The relationship between e_M and γ is explored in **Figure 10**. Elimination rate includes
 656 both data filtered out prior to AD as well as data labelled as anomalous by the AD
 657 method. A polynomial line of best fit has been added.

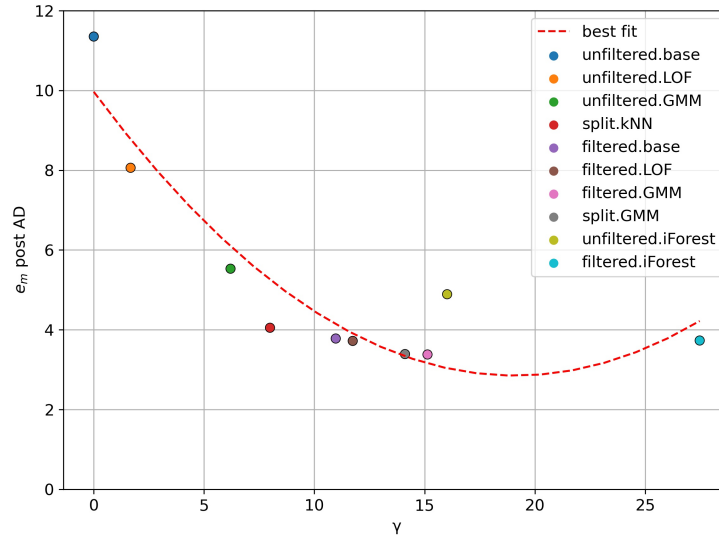


Figure 10: Results: Comparison of AD method averaged performance, γ versus e_m . On the x-axis are the percentages of SCADA data removed, γ , this includes both filtered and that eliminated by the AD technique. On the y-axis are prediction errors of the "cleaned" SCADA data, e_m , presented as a percentage of P_r .

658 Looking at these results, it is initially difficult to separate which is the greatest
 659 driver of reducing error, the approach or the method. From **Figure 10**, there appears
 660 to be a correlation between prediction error and amount of data removed, as shown by
 661 the line of best fit. This is intuitive, as removing anomalies should increase accuracy

662 and will reduce the amount of data remaining. However, the gradient of the line of
663 best fit shows that not all anomalies have equal influence on error. Using the line
664 of best fit, removing 10% of instances reduces error by 55.2% compared to retaining
665 all instances. Removing a further 10% of data only reduces error by a further 16.2%
666 (total reduction of 71.2%). Going beyond a γ of 17% appears to have no substantial
667 improvement on prediction error with the final method of *filtered.iForest* coming in at
668 a γ of 27.5% and an e_M of 3.73%. The high γ rates of the iForest methods appears to
669 be due to classifying large portions of the rated (flat) part of the WTPC as anomalous.
670 Note that this is a simple part of the curve to make power predictions for, hence the
671 generally low e_m is more impressive.

672 In terms of improvement to prediction error per percentage of data removed
673 ($\delta e_M/\gamma$), the *unfiltered.LOF* method’s performance is, by some margin, the best at
674 17.3 (see **Table 4**). This is slightly more than double that of the next best method,
675 which is *unfiltered.GMM* at 8.26. However, assessment of *unfiltered.LOF*’s perfor-
676 mance must be tempered by its poor absolute prediction error. Except for the base-
677 case itself (*unfiltered.base*), *unfiltered.LOF* ranks as the worst method for e_M . At a
678 high level, this trade-off suggests that different requirements of the SCADA data can
679 determine the choice in AD method. For example, if outright smallest prediction error
680 was preferred, no matter the elimination rate, *filtered.GMM* may be suitable. If the
681 preference was to retain data as far as possible then *unfiltered.LOF* may be selected.

682 4.3 Error and wind speed Variability

683 A similar trend can be seen for e_M and change to wind speed variability, δIQR_u , as
684 shown in **Figure 11**. This equates to an approximately linear relationship between γ
685 and δIQR_u .

686 From the results for *filtered.base* it appears that some reduction in IQR_u is accept-
687 able. This method reduces IQR_u by 5.5% and, looking at the WTPC in **Figure 9**, it
688 appears that this method has high precision but low recall. *Precision* is the measure
689 of how many instances labelled as ”anomaly” truly were anomalies. *Recall* is the total
690 amount of anomalies identified as a portion of all the anomalies i.e. true positives

691 versus the set of true positives and false negatives.

692 Comparing *filtered* and *unfiltered* methods, it is clear that the latter have the
 693 lowest δIQR_u . The *unfiltered* methods average a IQR_u of 5.66 (change of 3.5%),
 694 whilst *filtered* methods have an average of 5.33 (change of 9.1%). The *split* methods
 695 fair marginally better than *filtered* methods at an average of 5.39 (change of 8.14%).

696 Significantly further along the x-axis we have the two iForest methods. These
 697 methods have comparable prediction error results as others. Looking at their WTPCs
 698 in **Figure 9**, it seems the high δIQR for these SCADA sets is a result of labelling
 699 high and very low wind speeds as anomalous.

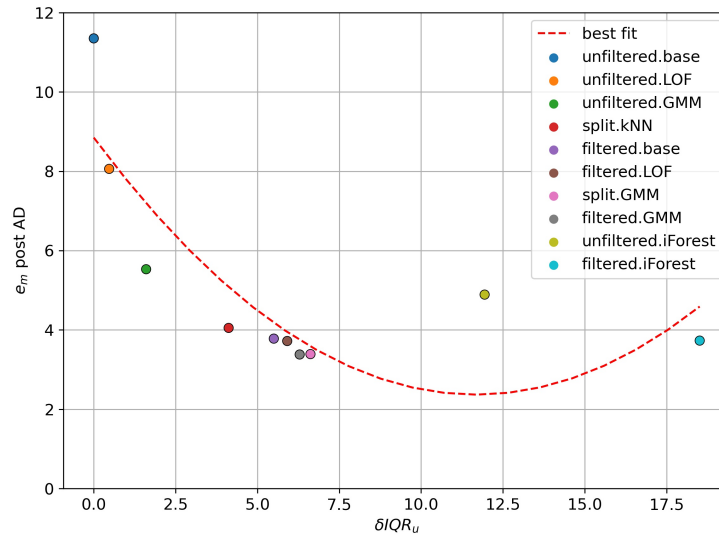


Figure 11: Results: Comparison of AD method averaged performance, δIQR_u versus e_m . On the x-axis are the percentage changes to windspeed IQR of each method, before and after it was applied, δIQR_u . On the y-axis are prediction errors of the "cleaned" SCADA data, e_m , presented as a percentage of P_r .

700 4.4 Advantages to Splitting?

701 As discussed previously, it is clear that from an absolute error perspective *filtered*
 702 produces far better results than *unfiltered*, but we must also consider *split*. As discussed
 703 in the methodology, the methods of *base*, *iForest*, and *LOF* are the same as for *split*

704 and *filtered* so cannot be used to isolate the impact of *split*. *GMM* differs between
705 *split* and *filtered*, and *kNN* is unique to *split*.

706 For *GMM*, *split* and *filtered* had broadly similar results. This is especially clear
707 when looking at **Figure 10**. Whilst the 2 approaches had similar results in terms of
708 error and IQR_u , *split* provided a very slight advantage with γ . *Filtered.GMM* removed
709 15.13% of data, marginally more than *split.GMM* with 14.11%. As such, *split.GMM*
710 may be appropriate in applications where high data retention is preferred.

711 For *split.kNN*, this method appeared to have good performance in terms of e_M/γ
712 and e_M/IQR_u but ultimately ranked poorly in terms of absolute e_M . A reason for
713 this can be found by comparing *split.kNN*'s initial class assignments (i.e. the same as
714 *filtered.base*) against final assignments, as shown in **Figure 12**. It appears the method
715 has incorrectly flipped many instances back from anomalous to normal so increasing
716 γ but reducing e_M . For the example WT shown in the figure, some 21,000 instances
717 were initially labelled as anomalous. Of these, approximately 9,000 (42%) flipped
718 class to "normal". Approximately 4,000 initially "normal" instances were assigned as
719 "anomalous". This is likely due to unbalanced class numbers whilst new instances
720 were being assigned. In the example, there were 3.5 times more initial "normal"
721 instances than "anomalies". This challenge should be accounted for in future research,
722 potentially through use of weightings.

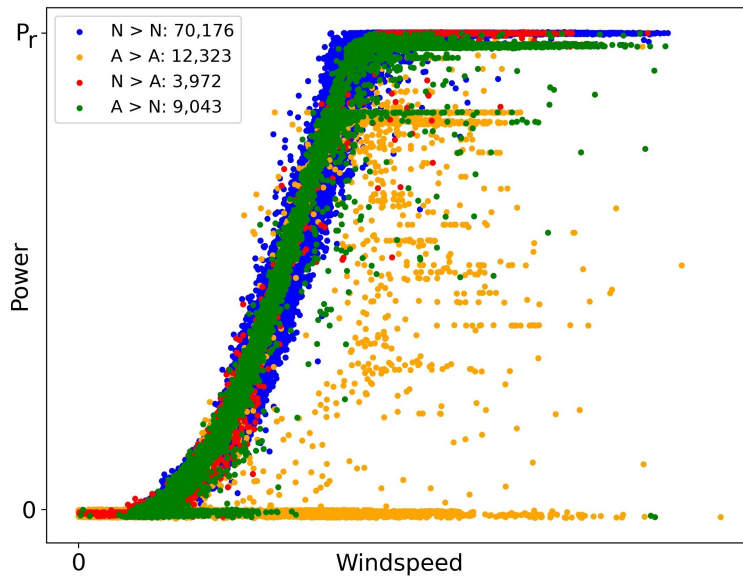


Figure 12: Visualisation of the change of classes from initial to final classification for the *split.kNN* method when applied to a WT from wind farm A. For the legend, "N" refers to "normal", "A" for "anomalous", and ">" the change from initial to final classification. The number of instances belonging to each category is shown in the legend.

723 4.5 Choice of AD Method

724 As discussed previously, the AD methods have been shown to have different perfor-
 725 mance characteristics in terms of prediction error, elimination rate, and change to the
 726 wind speed feature characteristics. The choice of AD method would therefore depend
 727 on the what the user of the SCADA data wishes to achieve. Two distinct scenarios
 728 are: (a) the lowest prediction error possible without unnecessary removal of SCADA
 729 data and change to wind speed variability; and (b) improving prediction error whilst
 730 maintaining as much SCADA data as possible.

731 For scenario (a), we start by examining *filtered.base*. As discussed previously, it
 732 appears likely that this method has high precision but low recall. As such, these results
 733 are used as a starting point from which to compare the other methods. This is shown
 734 in **Table 6**. Using this table, we can remove any methods with lower elimination rates
 735 on the ground of being less effective. These are: *unfiltered.base*, *unfiltered.GMM*, *unfil-*

Table 6: Results by method and relative to the *filtered.base* method. Note that *split* and *filtered* approaches are identical for the techniques of base, iForest, and LOF. Elimination rate includes both filtered data and that detected as anomalous.

Method	e_m	e_m relative to filtered.base	γ	γ relative to filtered.base	$IQR_{u,a}$	$IQR_{u,a}$ relative to filtered.base
unfiltered.base	11.35	-7.57	0.00	-10.97	5.86	0.32
unfiltered.iForest	4.89	-1.11	16.02	5.06	5.16	-0.38
unfiltered.GMM	5.53	-1.76	6.20	-4.77	5.77	0.23
unfiltered.LOF	8.06	-4.29	1.67	-9.30	5.84	0.30
filtered.base	3.78	-	10.97	-	5.54	-
filtered.iForest	3.73	0.05	27.47	16.50	4.78	-0.76
filtered.GMM	3.38	0.40	15.13	4.16	5.50	-0.05
filtered.LOF	3.72	0.05	11.74	0.77	5.52	-0.02
split.GMM	3.39	0.38	14.11	3.14	5.48	-0.07
split.kNN	4.05	-0.27	7.99	-2.98	5.62	0.08

736 *tered.LOF*, and *split.kNN*. Whilst this is not a guarantee that the remaining methods
737 eliminated the same instances as *filtered.base*, looking at **Figure 9**, this does appear to
738 be the case. Of the remaining 6 methods, the 2 iForest methods, (*filtered.iForest* and
739 *unfiltered.iForest*), can be removed as they have low rates of $\delta e_m / \delta IQR_u$ compared to
740 the other 4 methods, as per **Table 5**.

741 This leaves 4 unique methods in consideration. In order of e_M , these are: *fil-*
742 *tered.base* itself (3.78%), *filtered.LOF* (3.72%), *split.GMM* (3.39%), and *filtered.GMM*
743 (3.38%). The 2 GMM methods were clearly the more accurate and both are recom-
744 mended for scenario (a). As discussed in **Section 4.4**, it is slightly more advantageous
745 to use *split.GMM* over *filtered.GMM*; however, it should be noted that *split.GMM* was
746 considerably more challenging to implement than *filtered.GMM* due to the cross vali-
747 dations required (see **Section 2.3**). As such, some users may prefer *filtered.GMM*.

748 For scenario (b), we are concerned with methods with high rates of $\delta e_M / \delta IQR_u$
749 i.e. improvement to prediction error at minimal cost to the wind speed feature. The
750 3 methods with the highest rates in this category are: *unfiltered.LOF* (61.73), *unfil-*
751 *tered.GMM* (31.97), and *split.kNN* (15.61), according to **Table 5**. From **Figure 9**,
752 *unfiltered.LOF* was clearly overly conservative and can be eliminated. The remaining 2

753 methods showed similar results for their WTPCs; however, *unfiltered.GMM* is the bet-
754 ter choice with a δIQR_u approximately half that of *split.kNN* whilst still maintaining
755 a marked decrease in error from the base case (δe_M of 51.23%).

756 5 Conclusion

757 Three pre-processing approaches have been compared, along with 5 anomaly detection
758 techniques for a total of 10 unique AD methods. These methods have been applied to
759 SCADA data from 2 different wind farms for a total of 20 turbines. The efficacy of the
760 AD methods have been studied in terms of improvements to power prediction error,
761 amounts of data removed, and ability to retain the underlying statistical characteristics
762 of the wind speed feature. From this, and with respect to the SCADA sets used in the
763 study, the following conclusions are drawn:

- 764 • It is beneficial to pre-process the SCADA data by filtering out obvious anomalies
765 and explicit instances of faults/curtailments prior to applying anomaly detection
766 techniques.
- 767 • A pitch angle of $>30^\circ$ is proven to be a reasonable threshold for the above
768 pre-processing.
- 769 • All anomalies do not have equal impact upon error. The rate of prediction error
770 reduction reduces as more data is removed and anomalies become harder to
771 detect.
- 772 • The AD method of choice is dependent on the application, with some methods
773 achieving lower error at the cost of increasing percentages of data removal and
774 reduction in wind speed variability.
- 775 • The GMM technique is shown as an effective method to significantly reduce error
776 whilst maintaining statistical characteristics of wind speed data. This is espe-
777 cially so when combined with pre-processing anomalies, in which error reduces
778 by more than 70% compared to no pre-processing and no anomaly detection
779 technique.

- 780 • The *split.GMM* method appear to maintain marginally more data than its
781 *filtered* counterpart, however, the increased complexity in implementing this
782 method may make it undesirable.

783 Additionally, the importance of proper treatment of SCADA data regarding miss-
784 ing data has been raised. Given that SCADA data is the basis of so many findings,
785 conclusions, and concepts it is paramount that this treatment is discussed so that all
786 that follows can be relied upon or can be replicated by others.

787 6 Acknowledgments

788 The authors thank the anonymous data provider for providing the SCADA data for
789 both wind farms. This research was funded by the EPSRC Doctoral Training Part-
790 nership (EP/R513222/1 ENG).

791 References

- 792 [1] European Commission. *An EU Strategy to harness the potential of off-*
793 *shore renewable energy for a climate neutral future*. Tech. rep. European
794 Commission, 2020, p. 27. arXiv: arXiv:1011.1669v3. URL: https://ec.europa.eu/commission/presscorner/detail/en/IP_20_2096.
795
- 796 [2] UK Government. *Press Release: New plans to make UK world leader in*
797 *green energy*. Oct. 2019. URL: <https://www.gov.uk/government/news/new-plans-to-make-uk-world-leader-in-green-energy>.
798
- 799 [3] Varun Chandola, Arindam Banerjee, and Vipin Kumar. *Anomaly detec-*
800 *tion: A survey*. 2009. DOI: 10.1145/1541880.1541882.
- 801 [4] Zi Lin, Xiaolei Liu, and Maurizio Collu. “Wind power prediction based on
802 high-frequency SCADA data along with isolation forest and deep learning
803 neural networks”. In: *International Journal of Electrical Power & Energy*
804 *Systems* 118 (2020), p. 105835.

- 805 [5] Ran Bi, Chengke Zhou, and Donald M Hepburn. “Detection and classifica-
806 tion of faults in pitch-regulated wind turbine generators using normal be-
807 haviour models based on performance curves”. In: *Renewable Energy* 105
808 (2017), pp. 674–688. DOI: [https://doi.org/10.1016/j.renene.2016.](https://doi.org/10.1016/j.renene.2016.12.075)
809 [12.075](https://doi.org/10.1016/j.renene.2016.12.075). URL: [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/S0960148116311351)
810 [pii/S0960148116311351](https://www.sciencedirect.com/science/article/pii/S0960148116311351).
- 811 [6] Huan Long et al. “Image-based abnormal data detection and cleaning
812 algorithm via wind power curve”. In: *IEEE Transactions on Sustainable*
813 *Energy* 11.2 (2019), pp. 938–946.
- 814 [7] Xiaojun Shen, Xuejiao Fu, and Chongcheng Zhou. “A combined algorithm
815 for cleaning abnormal data of wind turbine power curve based on change
816 point grouping algorithm and quartile algorithm”. In: *IEEE Transactions*
817 *on Sustainable Energy* 10.1 (2018), pp. 46–54.
- 818 [8] Yue Wang et al. “Copula-based model for wind turbine power curve outlier
819 rejection”. In: *Wind Energy* 17.11 (2014), pp. 1677–1688.
- 820 [9] Le Zheng, Wei Hu, and Yong Min. “Raw wind data preprocessing: a
821 data-mining approach”. In: *IEEE Transactions on Sustainable Energy* 6.1
822 (2014), pp. 11–19.
- 823 [10] Yongning Zhao et al. “Data-driven correction approach to refine power
824 curve of wind farm under wind curtailment”. In: *IEEE Transactions on*
825 *Sustainable Energy* 9.1 (Jan. 2018), pp. 95–105. ISSN: 19493029. DOI: 10.
826 [1109/TSTE.2017.2717021](https://doi.org/10.1109/TSTE.2017.2717021).
- 827 [11] Neil Davis et al. “Ice detection on wind turbines using observed power
828 curve”. In: *Wind Energy* 19.6 (2016), pp. 999–1010. ISSN: 1095-4244. DOI:
829 [10.1002/we.1878](https://doi.org/10.1002/we.1878).

- 830 [12] IEC: TC 88 - Wind energy generation. *IEC TR 61400-12-4:2020*. Tech.
831 rep. IEC, 2020.
- 832 [13] Lijun Zhang et al. “Ice Detection Model of Wind Turbine Blades Based
833 on Random Forest Classifier”. In: *Energies* 11.10 (2018). ISSN: 1996-1073.
834 DOI: 10.3390/en11102548.
- 835 [14] Georgios Alexandros Skrimpas et al. “Detection of icing on wind turbine
836 blades by means of vibration and power curve analysis”. In: *Wind Energy*
837 19.10 (2016), pp. 1819–1832. DOI: <https://doi.org/10.1002/we.1952>.
- 838 [15] John W Graham, Patricio E Cumsille, and Allison E Shevock. “Methods
839 for handling missing data”. In: *Handbook of Psychology, Second Edition 2*
840 (2012).
- 841 [16] Fabian Pedregosa et al. “Scikit-learn: Machine learning in Python”. In:
842 *the Journal of machine Learning research* 12 (2011), pp. 2825–2830.
- 843 [17] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. “Isolation forest”. In:
844 *2008 eighth ieee international conference on data mining*. IEEE. 2008,
845 pp. 413–422.
- 846 [18] Brad Boehmke and Brandon M. Greenwell. *Hands-On Machine Learning*
847 *with R*. Boca Raton, FL, USA: CRC Press, 2019. ISBN: 978-1-138-49568-5.
- 848 [19] W. A. Shewhart. “Economic Quality Control of Manufactured Product”.
849 In: *Bell System Technical Journal* 9.2 (1930), pp. 364–389. ISSN: 15387305.
850 DOI: 10.1002/j.1538-7305.1930.tb00373.x.
- 851 [20] Jorma Laurikkala et al. “Informal identification of outliers in medical
852 data”. In: *Fifth international workshop on intelligent data analysis in*
853 *medicine and pharmacology*. 2000, pp. 20–24. URL: [https://scholar.
854 googleusercontent.com/scholar.bib?q=info:mZbjULJAXQoJ:scholar.
855 google.com/&output=citation&scisdr=CgXbp3N2ENqZ0IA7nlw:AAGBfm0AAAAAYHO-](https://scholar.googleusercontent.com/scholar.bib?q=info:mZbjULJAXQoJ:scholar.google.com/&output=citation&scisdr=CgXbp3N2ENqZ0IA7nlw:AAGBfm0AAAAAYHO-)

- 856 hlyr8srXhooBul_m-gqWthofg0u3&scisig=AAGBfm0AAAAAYH0-htKw6xhcq05qUs10F0AjRn10nEBc&
857 scisf=4&ct=citation&cd=-1.
- 858 [21] Aurélien Géron. *Hands-on machine learning with Scikit-Learn, Keras, and*
859 *TensorFlow: Concepts, tools, and techniques to build intelligent systems.*
860 O'Reilly Media, 2019.
- 861 [22] Markus M Breunig et al. “LOF: identifying density-based local outliers”.
862 In: *Proceedings of the 2000 ACM SIGMOD international conference on*
863 *Management of data.* 2000, pp. 93–104.
- 864 [23] E Fix and JL Hodges Jr. *Discriminatory Analysis. Nonparametric Dis-*
865 *crimination: Small Sample Performance.* Tech. rep. 1952. URL: <https://apps.dtic.mil/sti/citations/ADA800391>.
866
- 867 [24] Zi Lin and Xiaolei Liu. “Wind power forecasting of an offshore wind tur-
868 bine based on high-frequency SCADA data and deep learning neural net-
869 work”. In: *Energy* 201 (2020), p. 117693. ISSN: 0360-5442. DOI: <https://doi.org/10.1016/j.energy.2020.117693>. URL: <https://www.sciencedirect.com/science/article/pii/S0360544220308008>.
870
871