

Feature Extraction and Pattern Identification for Anemometer Condition Diagnosis

Longji Sun¹, Chao Chen², and Qi Cheng³

^{1,2,3} *School of Electrical & Computer Engineering, Oklahoma State University, Stillwater, OK, 74078, USA*

longji.sun@okstate.edu

chao.chen@okstate.edu

qi.cheng@okstate.edu

ABSTRACT

Cup anemometers are commonly used for wind speed measurement in the wind industry. Anemometer malfunctions lead to excessive errors in measurement and directly influence the wind energy development for a proposed wind farm site. In the PHM 2011 Data Challenge Competition, two types of data need to be processed for anemometer condition diagnosis: paired data consisting of wind data from paired anemometers, and shear data composed of measurements from an array of anemometers at different heights. Since the accuracy of anemometers can be severely affected by the environmental factors such as icing and the tubular tower itself, in order to distinguish the cause due to anemometer failures from these factors, our methodologies start with eliminating irregular data (outliers) under the influence of environmental factors. For paired data, the relation between the normalized wind speed difference and the wind direction is extracted as an important feature to reflect normal or abnormal behaviors of paired anemometers. Decisions regarding the condition of paired anemometers are made by comparing the features extracted from training and test data. For shear data, a power law model is fitted using the preprocessed and normalized data, and the sum of the squared residuals (SSR) is used to measure the health of an array of anemometers. Decisions are made by comparing the SSRs of training and test data. The performance of our proposed methods is evaluated through the competition website. As a final result, our team ranked the second place overall in both student and professional categories in this competition.

1. INTRODUCTION

Wind energy as a promising renewable energy source has attracted considerable attention in recent years. The first step in the development of a productive wind farm is wind re-

source assessment. Cup anemometers (IEA, 1999) have been widely used for wind speed measurement. Typical anemometers have three or four cups installed on a vertical shaft. Their measurements provide important information of wind resources for a proposed site. Therefore, their accuracy can greatly affect the estimated energy production and return on investment. Normally, the measurement of a cup anemometer is within 2% error. However, under some circumstances, such as the wear on the bearings, a missing cup or a failed shaft, an anemometer fails to provide accurate wind speed information, i.e., its measurements have excessive errors. It is critical that damaged or out of tolerance anemometers be detected and replaced in a timely manner.

Recent years have seen various methods proposed for the anemometer condition diagnosis problem. In (Beltran, Llombart, & Guerrero, 2009b), the nacelle anemometer fault detection problem is studied, in which wind speeds at one target anemometer are estimated by using two reference anemometers in its vicinity and the deviations of the estimates from the measurements are used to determine the target anemometer's condition. In (Beltran, Llombart, & Guerrero, 2009a), a method is introduced to select the range of data so that the uncertainty in evaluation of anemometers' health is minimized. To predict the failure of a hot-wire anemometer, a method utilizing a feature related sensor degradation and analyzing the trend of the feature is proposed (Delfino, Puttini, & Galvao, 2010). In the work by Kusiak, Zheng, and Zhang (2011), a virtual speed sensor is built based on historical wind speed data to monitor real sensors. In (Siegel & Lee, 2011), an anemometer assessment methodology using residual processing and clustering techniques is proposed, in which the residuals of anemometers' readings are computed and clustered to determine the anemometers' conditions.

The PHM 2011 Data Challenge is focused on the detection of failed anemometers. Generally, anemometers are installed on a meteorological tower. With single or paired anemometers at different heights, an array of anemometers is formed.

Longji Sun et.al. This is an open-access article distributed under the terms of the Creative Commons Attribution 3.0 United States License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

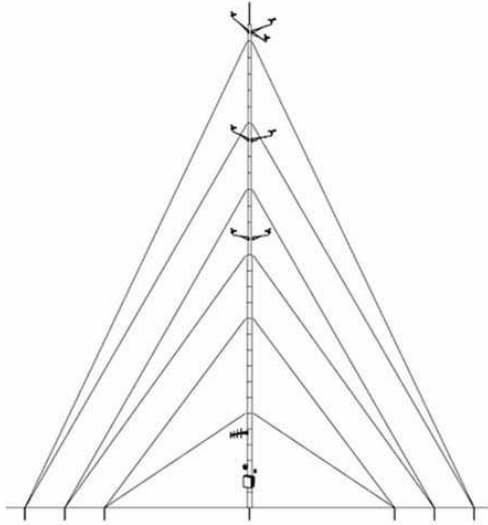


Figure 1. An example 60m tower with three sensor locations (<https://www.phmsociety.org/competition/phm/11/problem>).

Figure 1 shows an example of a 60m meteorological tower with sensors located at 59m, 49m and 39m, respectively. Figure 2 shows the normalized position of sensors at a tower. A paired data file includes the measurements of paired anemometers with a 90° or 180° angle, the corresponding wind direction and temperature. Each shear data file includes the measurements of an array of anemometers, wind direction, temperature and data collection time. The paired data consist of 12 training sets and the shear data consist of 7 training sets, each consists of 25 days of normal data. The paired data also have 420 test files and the shear data 255 test files, each of 5 days of data of unknown conditions. The problem is to detect failed anemometers in each test file. For paired data, it is to distinguish which one, or both of anemometers, fail if not both of them work normally. The objective for each shear test file is to determine whether all anemometers are in a good condition or not. Readers are referred to the PHM 2011 Data Challenge website (<https://www.phmsociety.org/competition/phm/11/problem>) for more information.

Since, in training files, only normal data are provided, the problem of anemometer fault detection is essentially anomaly detection. Various techniques have been developed for anomaly detection, including classification based methods (Duda, Hart, & Stork, 2000), statistical approaches (Barnett & Lewis, 1994), and clustering techniques (R. Smith, Bivens, Embrechts, Palagiri, & Szymanski, 2002). Anomaly can be categorized into point anomaly, contextual anomaly, and collective anomaly (Chandola, Banerjee, & Kumar, 2009). Point anomaly, i.e., anomalous individual data instance, is the most studied anomaly and the focus of most of the existing anomaly detection techniques. Contextual anomaly



Figure 2. Normalized positions of sensors on a tower (<https://www.phmsociety.org/competition/phm/11/problem>).

refers to a data instance only considered as an anomaly in a specific context. For example, in the work by Basu and Meckesheimer (2007), anomalies in time series data are detected by comparing the value of a data point with the median of its neighborhood. Collective anomaly means that a collection of data instances is anomalous, in which the relation between data is exploited to detect anomalies. For instance, sequential anomaly detection techniques are used to find unusual values in multiple time-series data (Chan & Mahoney, 2005).

In this paper, we will extract from training data important features that can reflect normal collective patterns or behaviors of anemometers in various contexts. Any deviation from these normal patterns can indicate possible faulty conditions. The rest of the paper is organized as follows. In Section 2, the methodology to analyze the paired dataset is provided. The method to deal with shear data is elaborated in Section 3. The paper is concluded with some discussion in Section 4.

2. METHODOLOGY FOR PAIRED DATA ANALYSIS

The method for paired data analysis mainly includes five steps: data preprocessing, feature extraction, denoising, pattern search and decision making. Firstly, a preprocessing step is taken to eliminate some apparently incorrect and statistically useless measurements. Secondly, a feature, namely, the relation between the discrepancy of the paired anemometer measurements and the wind direction, is extracted from the preprocessed data. A further denoising step is taken to reduce the environmental effects and make the feature more prominent in different situations. Then, an algorithm is designed to search for each test data file the most matched pattern from training data. Finally, decisions are made based on the rela-

tion between the pattern under testing and the matched pattern.

2.1. Data Preprocessing

Failed anemometers cannot provide accurate wind speed measurements. On the other hand, environmental factors, such as icing can also affect the accuracy of measurements considerably. To avoid false alarms, it is important to distinguish these two types of situations. Some preprocessing of the raw measurements is required.

The preprocessing step is composed of two stages. In the first stage, data undergo a measurement range test. Namely, only measurements within a valid measurement range are meaningful. Factors, such as sensor noise and icing, result in measurements outside this range, which fail to provide useful information and should be eliminated. For this problem, the range is set to be from 0.4m/s to 75m/s (<https://www.phmsociety.org/competition/phm/11/problem>).

In the second stage, detection of icing conditions is conducted. Icing is a leading factor in introducing errors in measurement data. Empirical results (Kenyon & Blittersdorf, 1996) and our observations of the training data have shown that icing conditions have the following characteristics:

- 1) When the temperature is at or below the icing point, the standard deviation of the wind speed measurements is zero or near zero.
- 2) The standard deviation of the wind direction measurements is zero or near zero.

In (Schaffner, 2002), it is suggested that the measurements in six hours before and after the icing points should be discarded, considering that the effect of icing begins long before an anemometer is frozen and continues for some time before the frozen effect completely disappears. Since we have limited data in this competition, especially for test data, a more practical range is adopted in which only the data in 30 minutes before and after icing points are discarded.

2.2. Feature Extraction

In an ideal environment, the measurements of a pair of normal anemometers should be very close to each other given that they measure the wind speeds at the same height with a very close distance. However, this is not always the case for the given training data. It can be shown that the mast of the tubular tower on which the paired anemometers are mounted plays an important role (Lubitz, 2009). The mast of the tower will generate a wake behind it, acceleration around it and a retardation upwind of it (IEA, 1999). Figure 3 shows a wind field around the mast of a tubular tower. The numbers indicate the ratio between local wind speeds and the free-field wind speed. This fact explains the significant difference in

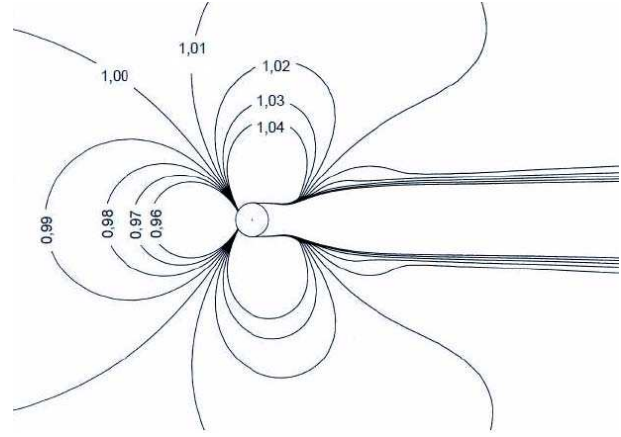


Figure 3. The wind field around a tubular mast (from the IEA 1999 report).

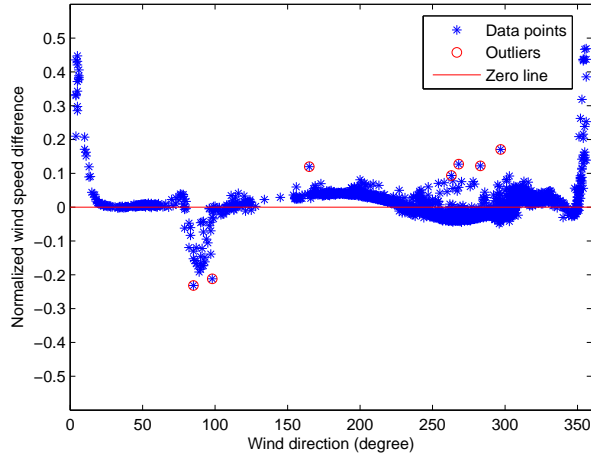
paired anemometer measurements in some wind directions. This suggests that the relation between the wind speed difference and its corresponding wind direction can be utilized as a key feature to describe the condition of paired anemometers. The wind speed difference is computed as follows:

$$s = \frac{s^{(1)} - s^{(2)}}{\max(s^{(1)}, s^{(2)})} \quad (1)$$

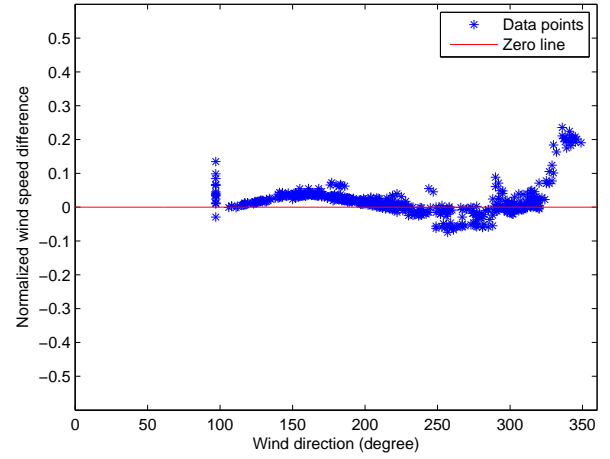
where $s^{(1)}$ is the wind speed of anemometer 1 and $s^{(2)}$ is the wind speed of anemometer 2. Normalization is taken to simplify the subsequent pattern search step. This is different from (Lubitz, 2009), where $s = s^{(1)}/s^{(2)}$ is used as a wind difference indicator to evaluate the tower effect approximation model. Figures 4(a) and 4(b) show the normalized wind speed difference as a function of the wind direction for pairTrng1 and pairTrng7 for example. Figures 5(a) and 5(b) plot the same relation for two test data files. Since the training data are from normal anemometers, the relations between the wind speed difference and the wind direction based on these training files are the representatives of normal behaviors of anemometers. Deviations from these representative patterns may indicate failure of anemometers in test data.

2.3. Denoising

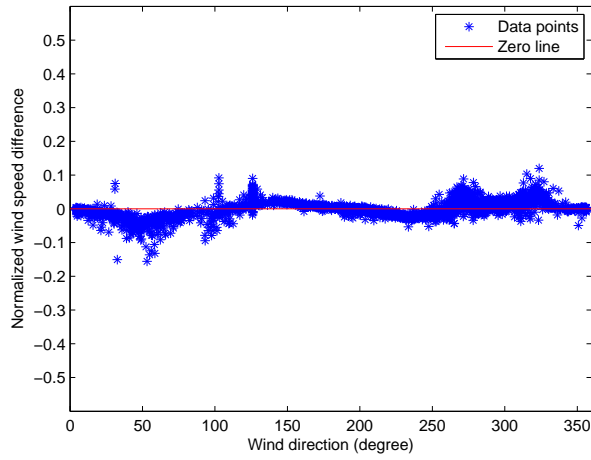
In Figure 4(a), we observe that around 90° and 360° , the wind speed difference deviates from zero, while for the rest of wind directions, the difference varies around zero. This may be due to the normalized position of the paired anemometers with respect to the mast. Besides, there are some data points isolated from the majority of the rest, which are marked with circle in the figure. This situation is more severe in test data. Because of the limited size, the percentage of isolated points in test data can be large. By checking the original data, the isolated data points generally correspond to a low wind temperature when anemometers may run slow. This is the case for all training data. To make the pattern more prominent and make



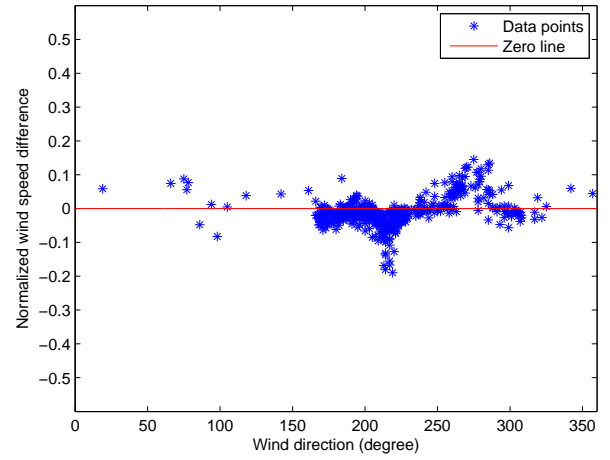
(a) PairTrng1



(a) Pairdata1



(b) PairTrng7



(b) Pairdata12

Figure 4. Normalized wind speed difference as a function of wind direction for training data.

it easier for test data to find a match, we consider these isolated points as outliers and they should be eliminated.

There are many ways to remove outliers. The method adopted here is based on the average distance of each data point from its k nearest neighbors, the larger of which indicates it is more likely to be an outlier. More information about the distance-based outlier detection techniques can be found in (Knorr, Ng, & Tucakov, 2000). The distance between two data points in Figure 4(a) is defined as follows:

$$D = \sqrt{\left(\frac{d_i - d_j}{360}\right)^2 + (s_i - s_j)^2} \quad (2)$$

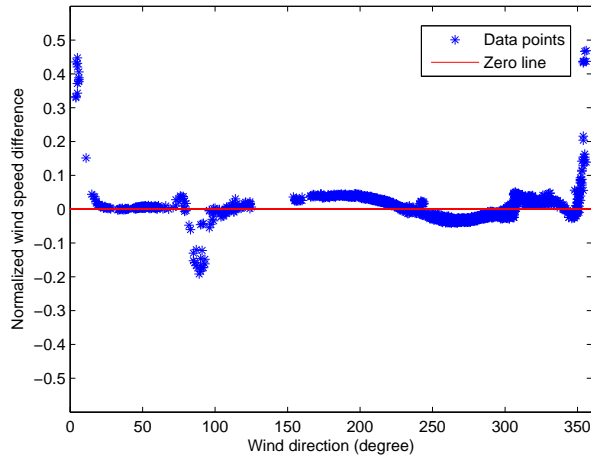
where d_i (d_j) is the wind direction and s_i (s_j) is the wind speed difference along that direction. Normalizing d_i by 360° is to make these two quantities comparable. For every data

Figure 5. Normalized wind speed difference as a function of wind direction for test data.

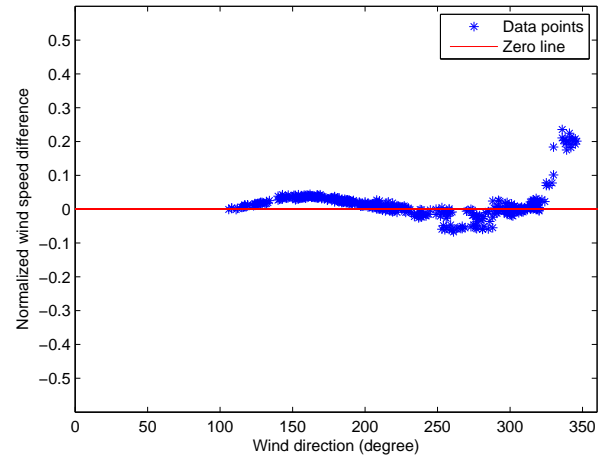
point (d, s) , the average distance from its k nearest neighbors is calculated,

$$D_{(d,s)} = \frac{1}{k} \sum_{i=1}^k \sqrt{\left(\frac{d - d_i}{360}\right)^2 + (s - s_i)^2} \quad (3)$$

where $\{(d_i, s_i), i = 1, \dots, k\}$ is the set of k nearest neighbors. Since the distribution of data points is different at different wind direction, we compare $D_{(d,s)}$ only with that of those data points of similar wind directions. A window of length Δd moves along the wind direction axis. For all data points in this window, those whose average distance is among the largest $\alpha\%$ are marked as outliers and are eliminated. The performance of this method depends on parameter k , Δd and α . In our experiments, we set $k = 10$, $\Delta d = 20$ and $\alpha = 10$ which gives good empirical results. Figures 6(a) and 6(b)



(a) PairTrng1



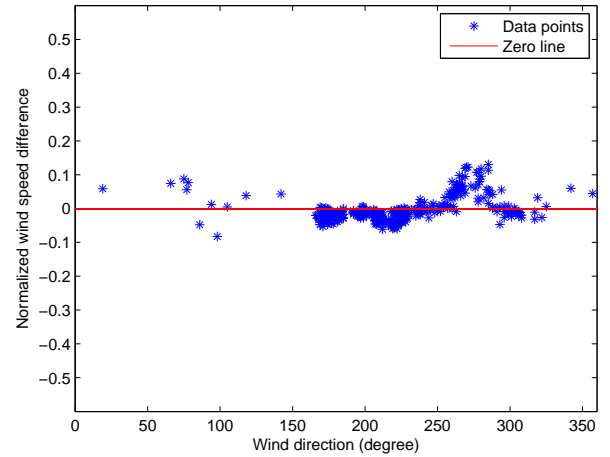
(b) PairTrng7

Figure 6. Normalized wind speed difference vs wind direction for training data after denoising.

are the relation of wind speed difference and wind direction after removing outliers (isolated points). This can also be applied to test data files and the results are shown in Figures 7(a) and 7(b).

2.4. Pattern Search

Training data are collected from normal anemometers. Since there are twelve training files, there are twelve normal patterns under different configurations. In this step, we need to find, for each test file, the most matched training profile for comparison. Distance is the most used metric to measure the similarity of two patterns. Given a training file p and a test file q , assume that there are a total of N wind directions that both training and test files have wind speed difference values. If these values are plotted in an N -dimensional space,



(a) Pairdata1

Figure 7. Normalized wind speed difference vs wind direction for test data after denoising.

two point clouds are formed. Figure 8 is an example when $N = 2$. The distance of two point clouds can be measured by the distance between their centroids. More specifically,

$$Dis(q, p) = \frac{1}{N} \|\bar{S}_q - \bar{S}_p\|_2 \quad (4)$$

where \bar{S}_q is an N -dimensional vector, each element of which is the mean wind speed difference for that wind direction. The same applies for \bar{S}_p , $p = 1, \dots, 12$. Normalization over N is done to eliminate the effect of the number of dimensions. Another important factor is the shape of data distribution pattern. The similar shape indicates a similar anemometer configuration. The correlation coefficient is adopted and defined

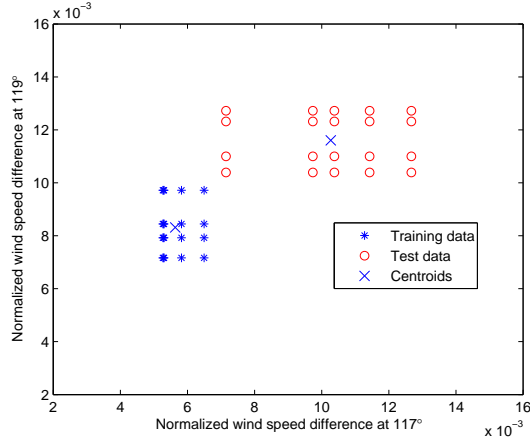


Figure 8. Two point clouds for training data (*) and test data (o) when wind directions 117° and 119° are selected. “x” represent the centroids of the clouds.

as follows:

$$\rho(q, p) = \frac{\langle \bar{S}_q, \bar{S}_p \rangle}{\|\bar{S}_q\|_2 \|\bar{S}_p\|_2} \quad (5)$$

Here, $\langle \cdot, \cdot \rangle$ stands for the inner product of two vectors. The larger the correlation coefficient, the more similar the two shapes of p and q are. The training profile p^* is selected for comparison with test file q if

$$p^* = \arg \min_p \left(\frac{Dis(q, p)}{\|\bar{S}_q\|_2} + \sqrt{1 - \rho^2(q, p)} \right) \quad (6)$$

The objective function is the average of the distance measure and shape measure.

2.5. Decision Making

There are four possible conditions of the paired anemometers in the test data: both are normal (0), anemometer 1 fails (1), anemometer 2 fails (2), and both fail (3). Following assumptions are made regarding these four conditions:

- (0) If both anemometers work normally, the feature, i.e., the relation between the wind speed difference and the wind direction, should be very similar to its corresponding matched training pattern. That is, the feature extracted from the test data file will have a significant overlap with the corresponding training data pattern. Figures 9¹ is an example.
- (1) It is assumed that if an anemometer fails, its reading is generally smaller than the true value, especially for mechanical failures. Based on the definition of the wind speed difference in Eq. (1), if anemometer 1 fails, the pattern will have a downward shift. Namely, the wind speed difference values will take more negative values with the change of wind direction.

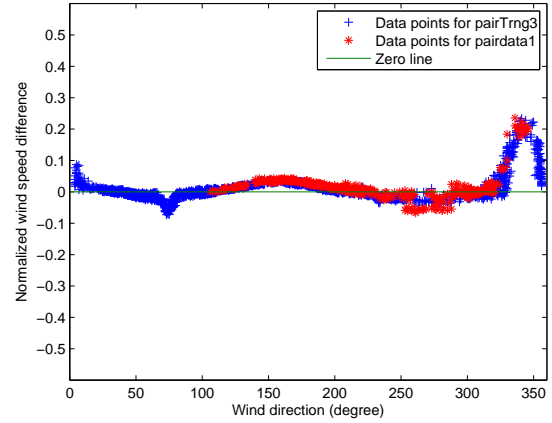


Figure 9. Significant overlap of patterns extracted from pairdata1 and pairTrng3.

- (2) With the same assumption, if anemometer 2 fails, the pattern shows an upper shift. That is, the wind speed difference values will take more positive values. There are many such kinds of patterns in test data, which shows that this assumption may be right. Figures 10 and 11 are the examples of these two conditions.
- (3) If both anemometers fail, the pattern is not predictable, i.e., it does not show any of the above characteristics in an obvious way.

To make a decision for each test file, the following algorithm is designed taking into account the above assumptions. Assume that for test file q , training file p is selected for comparison through the pattern search step. For wind direction d where wind speed difference values are available in both training and test data, define $S_p(d)_{min} = \min\{S_{p_1}(d), \dots, S_{p_n}(d)\}$ and $S_p(d)_{max} = \max\{S_{p_1}(d), \dots, S_{p_n}(d)\}$, assuming that there are n wind speed difference values at wind direction d in training file p . Then for test data q and for the same direction d , count the number of data points in, above or below the range $[S_p(d)_{min}, S_p(d)_{max}]$, which are denoted as $C_{q,in}(d)$, $C_{q,above}(d)$, and $C_{q,below}(d)$, respectively. There are two ways to proceed based on these counts. One is to make a decision for each wind direction and fuse these decisions to generate a global decision (decision fusion). The other one is to add up the total number of data points in, above or below the normal ranges and make a decision based on that (data fusion). Since we have no ground truth and the characteristics of wind speed difference vary for different wind directions, we develop the following hybrid method. The whole 360° is divided into 36 bins. The counts in each bin add up, i.e.,

$$C_{q,xx}(i) = \sum_{d \in Bin_i} C_{q,xx}(d) \quad (7)$$

¹Figures 9, 10, and 11 can be viewed better with a color print.

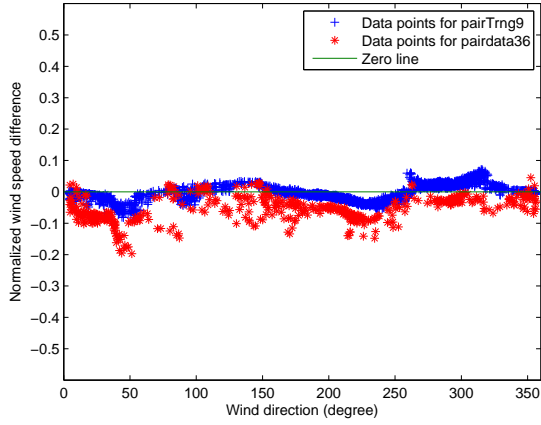


Figure 10. Significant downward shift of the pattern extracted from pairdata36 compared to the pattern from pairTrng9.

where $i = 1, \dots, 36$ and xx can be *in*, *above* and *below*. Decision is made for each bin using the following rule:

$$U_q(i) = \begin{cases} 0 & \text{if } C_{q,in}(i) > T_i \\ 1 & \text{if } C_{q,below}(i) > T_i \\ 2 & \text{if } C_{q,above}(i) > T_i \\ 3 & \text{otherwise} \end{cases} \quad (8)$$

where threshold $T_i = \frac{C_{q,in}(i) + C_{q,above}(i) + C_{q,below}(i)}{2}$. That is, whichever of the first three conditions dominating indicates the condition of that bin. If there is no one that dominates, decision 3 is made. The majority of local decisions is chosen as the global decision. This hybrid method can not only smooth out the noise effect, but also preserve the variation of data pattern in different directions. Note that if no data points exist in some bins, those bins do not participate in decision making.

2.6. Results and Discussion

In the competition, the results are evaluated based on whether the proposed algorithm can accurately determines the conditions of the paired anemometers for each test file. Credit for each file is gained only if the decisions for both anemometers are correct. Visualization of our results for paired data is provided in Figure 12 on the top of next page. Condition indicators 0, 1, 2, and 3 are defined in Section 2.5. There are a total of 287 test files with decision 0, 43 files with decision 1, 39 files with decision 2, and 51 files with decision 3.

For paired data analysis, the normalized wind speed difference (NWS) as a function of the wind direction is extracted as a main feature for the purpose of faulty anemometer detection. Since wind data are collected from different environments, under different weather conditions and with different tower configurations, taking the difference and normalization of paired data can reduce environmental impacts effectively,

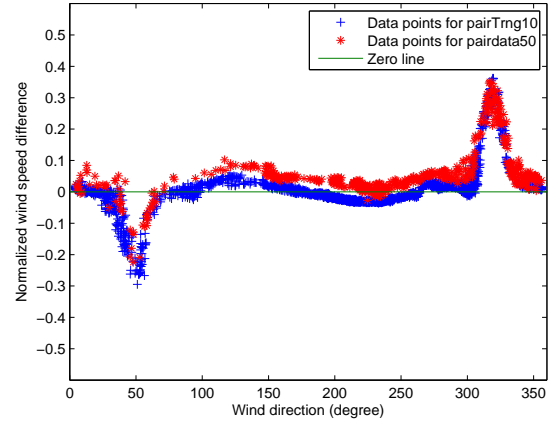


Figure 11. Upper shift of the pattern extracted from pairdata50 compared to the pattern from pairTrng10.

while the NWS pattern with respect to the wind direction can help identify similar anemometer configurations, thus putting training/testing-file comparison and anomaly detection in the same context. If raw data instead of the proposed feature is used, almost all test files look different/erroneous compared to the training files.

3. METHODOLOGY FOR SHEAR DATA ANALYSIS

For shear data, the problem is to decide whether all of an array of anemometers work normally. Similarly, a data preprocessing step has to be taken to eliminate some obviously useless data. Specifically, the measurement range test is conducted. It should be noted that the effect of icing conditions in cold climate is huge so that a majority of data are under the influence to different extents (Schaffner, 2002). For this problem, the same criteria as specified for paired data are used to partially mitigate the icing effect.

3.1. Irregular Data Elimination

Generally, the wind speed increases with the height because of the wind shear effect. However, in the training data with all anemometers in a normal condition, there exist many measurements violating this rule. This indicates that the measurements do not always reflect the true wind speeds, which may be due to the environmental factors rather than anemometer failures. We define a record containing this kind of measurements as irregular data and they make the detection problem more challenging. In Table 1, we summarize the mean temperature and the percentage of irregular data for all 7 shear training files. It is noted that the ones with lower temperatures generally have more irregular data. Thus, the irregularity is more likely the result of icing effects. To reduce the effect of icing on decision making, we eliminate all the irregular data.

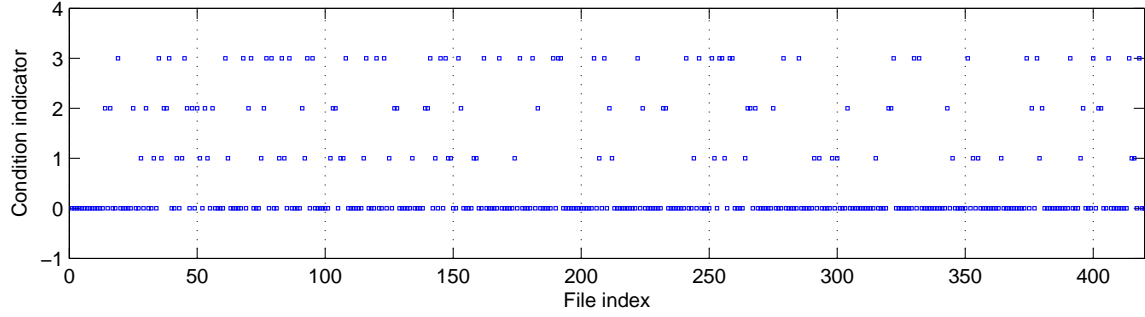


Figure 12. Results for paired data.

ShearTrng file	Mean Temperature (°F)	Irregular Data (%)
1	11.45	11
2	50.81	8
3	45.20	3
4	3.07	67
5	3.23	68
6	11.02	35
7	10.93	35

Table 1. Mean temperature and the percentage of irregular data for shear training files.

As mentioned in Section 2, the configuration of the tower also has effects on wind speed measurements. In Figure 13, the ratio of wind speeds at 59m and 51m as a function of the wind direction is plotted. It shows that around 90°, the direction in which the anemometers are installed, the ratio takes significantly different values. The measurements around the wind directions, to which the anemometers are pointed, fail to reflect the normal situation and therefore are eliminated.

3.2. Model fitting

After eliminating irregular data due to icing and/or the tower, we assume the failure of anemometers is the dominating factor of irregular patterns in test data, if any. One widely used wind shear model is a power law model (Burton, Sharpe, Jenkins, & Bossanyi, 2001), and is given as follows,

$$\frac{s}{s_r} = \left(\frac{h}{h_r} \right)^\alpha \quad (9)$$

where s is the wind speed at some specific height h , s_r the wind speed at a reference height h_r , and α the shear exponent. If we use the shear data to fit the model, we expect that the sum of squared residuals (SSR) tends to be small for normal data while be relatively large for abnormal data. Since the wind speed changes across time and space, to make the SSR

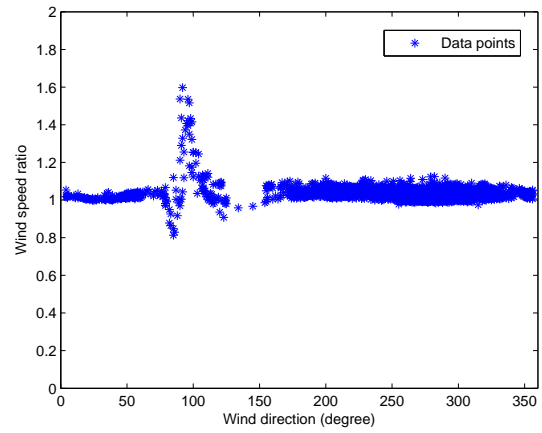


Figure 13. Wind speed ratio between 59m and 51m versus wind direction for shearTrng1.

comparable for different files, normalization by the maximum value of an array of wind speeds at each time is taken. As a result, all normalized wind speed fall into the range of $[0, 1]$. Figure 14 shows an example of the fitted power law model and normalized sample data points for shearTrng1. SSR is used as a performance measure of the given shear data.

3.3. Decision Making

There are three types of shear data files: three anemometers at (57m, 45m, 35m), four anemometers at (59m, 51m, 30m, 10m) and four anemometers at (49m, 39m, 30m, 10m). For the first two types, the training and test data have very similar temperature. Since there is only one and two training files for these two types of data respectively, a 25-day training data file is divided into 5 files with 5 days of data each. The SSRs are calculated as shown in Figure 15 for the four-anemometer configuration for 5 smaller training files and 20 test files. The decision making rule is as follows. The average of five SSR values from training data is used as a threshold, partially eliminating the randomness such as noise. If

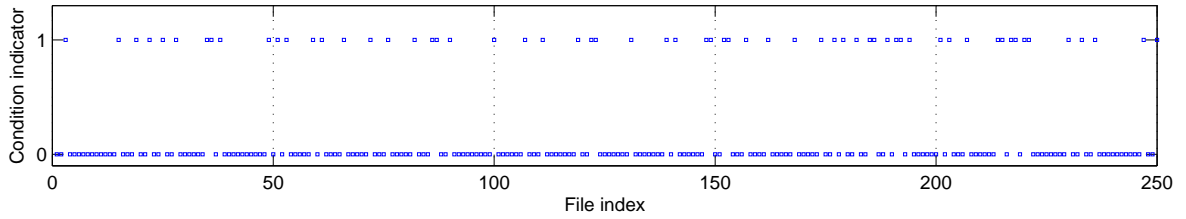


Figure 16. Results for shear data

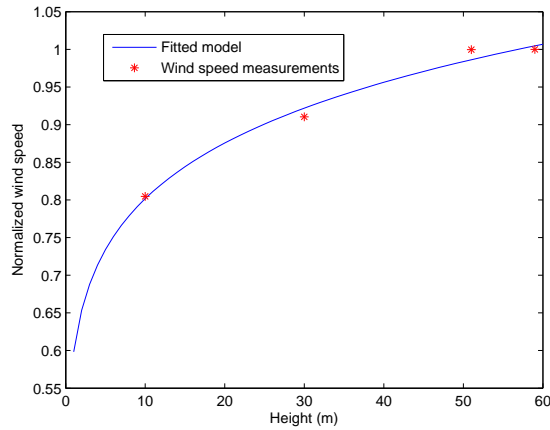


Figure 14. Normalized wind speeds measurements and fitted wind speeds using a power law model.

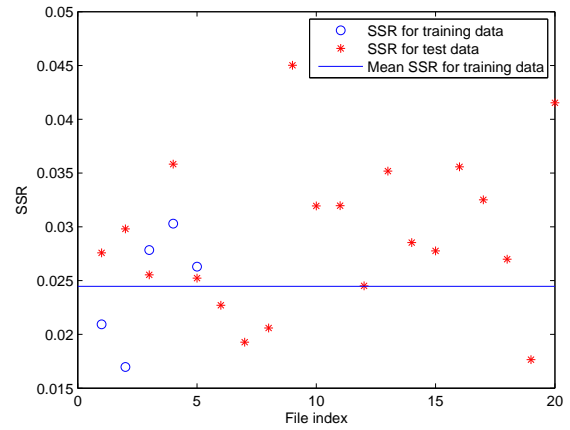


Figure 15. SSR for training and testing data for the four-anemometer configuration (59m, 51m, 30m, 10m).

the SSR of a test file is greater than this threshold, the decision that not all anemometer working normally is made. For the configuration with (49m, 39m, 30m, 10m), there are two types of temperature values: below the icing points and much above the icing points. This fact motivates us to compare the test data with the training files of similar temperature. The rest of algorithm remains the same.

3.4. Results and Discussion

In the competition, the results are evaluated based on whether the proposed algorithm can accurately determine the condition of an array of anemometers for each test file. Credit for each file is gained if the decision about whether any faulty anemometer occurs among the array of anemometers is correct. Visualization of our results for shear data is provided in Figure 16. A test file without any faulty anemometer is indicated with number 0, otherwise with number 1. There are a total of 193 files with decision 0 and 62 files with decision 1.

For shear data, the performance of the proposed algorithm largely depends on the elimination of noisy and irregular data. The sum of squared residuals (SSR) after fitting a power law model for an array of normalized wind speed measurements is used as the main feature to detect if any fault exists in the

array. The assumption is that a faulty array tends to have a larger SSR compared to that of training data of the same configuration. Data of different time of a day used for model fitting may influence the decision. This is because the wind shear is also a function of the time of a day and exhibits the diurnal variation, i.e., the wind shear exponent in the daytime is significantly smaller than at night (K. Smith, Randall, Malcolm, Kelley, & Smith, 2002). Therefore, an improved feature may be SSR as a function of the time of a day, the pattern variation of which can also be used as an indicator of possible faults. This will be investigated in our future work.

4. CONCLUSIONS

In this paper, we have developed a series of methods including data preprocessing, feature extraction and pattern identification to solve the anemometer condition diagnosis problem of the PHM 2011 Data Challenge Competition. The main idea of the algorithms is to extract useful features showing discernable patterns of training and test data so that they can reflect the health condition of anemometers. Since the data patterns may also be significantly influenced by various factors such as icing and the tower rather than anemometer failures, considerable efforts have been taken for elimination of

irregular data due to these environmental factors. For paired data, the relation between the normalized wind speed difference and the wind direction is used as the key feature for pattern identification and decision making. For shear data, the sum of squared residuals after model fitting is used for decision making. Several important assumptions are made for algorithm development, some of which have been justified by our observation of the data and the domain knowledge.

There are several aspects that we can pursue to further improve the diagnosis performance: the development of more efficient methods to reduce environmental effects and eliminate outliers, e.g., a new criterion by looking at wind direction to determine the range of useless data; extraction of features that are more sensitive to anemometer failures. The influence from environment is a major challenge which prevents us from accurately capturing the characteristics of anemometer failures. On the other hand, features more sensitive to anemometer failures would lead to a higher failure detection probability and a lower false alarm rate. These efforts will all have great practical values to the wind energy development.

ACKNOWLEDGMENT

This work was supported by NSF under Grant CPS-0932297.

REFERENCES

- Barnett, V., & Lewis, T. (1994). *Outliers in Statistical Data (3rd Edition)*. John Wiley and Sons.
- Basu, S., & Meckesheimer, M. (2007). Automatic Outlier Detection for Time Series: An Application to Sensor Data. *Knowledge and Information Systems*, 11(2), 137–154.
- Beltran, J., Llombart, A., & Guerrero, J. J. (2009a). A Bin Method with Data Range Selection for Detection of Nacelle Anemometers faults. In *Proceedings of European Wind Energy Conference and Exhibition (EWEC)*. March 16-19, Marseille, France,.
- Beltran, J., Llombart, A., & Guerrero, J. J. (2009b). Detection of Nacelle Anemometers Faults in a Wind Farm. In *Proceedings of International Conference on Renewable Energies and Power Quality (ICREPQ)*. April 15-17, Valencia, Spain.
- Burton, T., Sharpe, D., Jenkins, N., & Bossanyi, E. (2001). *Wind Energy Handbook (2nd Edition)*. Wiley.
- Chan, P., & Mahoney, M. (2005). Modeling Multiple Time Series for Anomaly Detection. In *Proceedings of Fifth IEEE International Conference on Data Mining* (pp. 90–97). November 27-30, Houston, TX, USA.
- Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly Detection: a Survey. *ACM Computing Surveys*, 41(3), 1–58.
- Delfino, T. N., Puttini, L. C., & Galvao, R. K. H. (2010). Fault Prognosis of an Air Flow Sensor. In *Proceedings of XVIII Congresso Brasileiro de Automtica (CBA)*. September 12-16. Bonito, MS, Brazil.
- Duda, R. O., Hart, P. E., & Stork, D. G. (2000). *Pattern Classification (2nd Edition)*. Wiley-Interscience.
- IEA. (1999). *Annex XI: Recommended Practices for Wind Turbine Testing and Evaluation 11. Wind Speed Measurement and Use of Cup Anemometry, I*. Paris: IEA.
- Kenyon, P. R., & Blittersdorf, D. C. (1996). *Accurate Wind Measurements in Icing Environments, Solutions to the Problem of Invalid Data from Frozen Anemometers and Direction Vanes*. Report NRG System.
- Knorr, E. M., Ng, R. T., & Tucakov, V. (2000). Distance-based Outliers: Algorithms and Applications. *The Very Large Data Bases (VLDB) Journal*, 8(3-4), 237–253.
- Kusiak, A., Zheng, H., & Zhang, Z. (2011). Virtual Wind Speed Sensor for Wind Turbines. *Journal of Energy Engineering*, 137(2), 59–69.
- Lubitz, W. D. (2009). Effects of Tower Shadowing on Anemometer Data. In *Proceedings of 11th Americas Conference on Wind Engineering*. June 22-26, San Juan, Puerto Rico.
- Schaffner, B. (2002). *Wind Energy Site Assessment in Harsh Climatic Conditions, Long Term Experience in Swiss Alps*. Report METEOTEST.
- Siegel, D., & Lee, J. (2011). An Auto-Associative Residual Processing and K-means Clustering Approach for Anemometer Health Assessment. *International Journal of Prognostics and Health Management*, 2(2).
- Smith, K., Randall, G., Malcolm, D., Kelley, N., & Smith, B. (2002). Evaluation of Wind Shear Patterns at Midwest Wind Energy Facilities. In *Proceedings of the American Wind Energy Association (AWEA) Windpower 2002 Conference*. June, Portland, OR, USA,.
- Smith, R., Bivens, A., Embrechts, M., Palagiri, C., & Szymanski, B. (2002). Clustering Approaches for Anomaly Based Intrusion Detection. In *Proceedings of Intelligent Engineering Systems through Artificial Neural Networks*. (pp. 579–584). ASME Press.
- Longji Sun** received his B.E. degree in communication engineering from University of Shanghai for Science and Technology, Shanghai, China in 2010. He is now pursuing his Ph.D. degree in the School of Electrical and Computer Engineering at Oklahoma State University. His current research interests include structural health monitoring using wireless sensor networks, cooperative optimization and cognitive radios. He is a student member of IEEE.
- Chao Chen** received his B.E. degree in automation control systems from Wuhan University of Technology, Wuhan, China in 2007. From 2007 to 2008, he worked as an electrical system design engineer at Shanghai Waigaoqiao Shipbuilding Company, Ltd., Shanghai, China. He is currently an M.S student in the School of Electrical and Computer Engineering at Oklahoma State University. His research interests include

structural health monitoring, signal processing, information fusion and statistical analysis.

Qi Cheng received the B.E. degree in electrical engineering (highest honors) from Shanghai Jiao Tong University, Shanghai, China, in July 1999, and the M.S. and Ph.D. degrees in electrical engineering from Syracuse University, Syracuse, NY, in 2003 and 2006, respectively. From 1999 to 2000, she worked as a System Engineer at Guoxin Lucent Technolo-

gies Network Technologies Company, Ltd., Shanghai, China. Since August 2006, she has been with Oklahoma State University, as an Assistant Professor in the School of Electrical and Computer Engineering. Her area of interest mainly focuses on statistical signal processing and data fusion with applications in wireless communications and distributed sensor networks. She is a member of IEEE.