

Article

SCADA Data-Based Support Vector Machine Wind Turbine Power Curve Uncertainty Estimation and Its Comparative Studies

Ravi Pandit ^{1,*}  and Athanasios Kolios ²

¹ Computer Science Department, University of Exeter, Exeter EX4 4PY, UK

² Naval Architecture, Ocean & Marine Engineering Department, University of Strathclyde, Glasgow G1 1XQ, UK; athanasios.kolios@strath.ac.uk

* Correspondence: r.pandit@exeter.ac.uk

Received: 6 November 2020; Accepted: 3 December 2020; Published: 4 December 2020



Featured Application: Proposed research could be useful in offshore wind turbine condition monitoring activities based on supervisory control and data acquisition (SCADA) datasets.

Abstract: Power curves, supplied by turbine manufacturers, are extensively used in condition monitoring, energy estimation, and improving operational efficiency. However, there is substantial uncertainty linked to power curve measurements as they usually take place only at hub height. Data-driven model accuracy is significantly affected by uncertainty. Therefore, an accurate estimation of uncertainty gives the confidence to wind farm operators for improving performance/condition monitoring and energy forecasting activities that are based on data-driven methods. The support vector machine (SVM) is a data-driven, machine learning approach, widely used in solving problems related to classification and regression. The uncertainty associated with models is quantified using confidence intervals (CIs), which are themselves estimated. This study proposes two approaches, namely, pointwise CIs and simultaneous CIs, to measure the uncertainty associated with an SVM-based power curve model. A radial basis function is taken as the kernel function to improve the accuracy of the SVM models. The proposed techniques are then verified by extensive 10 min average supervisory control and data acquisition (SCADA) data, obtained from pitch-controlled wind turbines. The results suggest that both proposed techniques are effective in measuring SVM power curve uncertainty, out of which, pointwise CIs are found to be the most accurate because they produce relatively smaller CIs. Thus, pointwise CIs have better ability to reject faulty data if fault detection algorithms were constructed based on SVM power curve and pointwise CIs. The full paper will explain the merits and demerits of the proposed research in detail and lay out a foundation regarding how this can be used for offshore wind turbine conditions and/or performance monitoring activities.

Keywords: wind turbines; power curves; SCADA datasets; condition monitoring; machine learning; support vector machines

1. Introduction

Wind energy in recent years has gained popularity because of low life cycle emissions and efforts to reduce costs. From a business perspective, the cumulative installed wind energy capacity globally is anticipated to reach 817 GW in 2021, where Asia will be leading with an installed capacity of 153.5 GW in the years 2017 to 2021, as reported by the Global Wind Energy Council (GWEC) [1]. Due to technology maturity, an exponential increase in wind turbine (WT) installation has been recorded. Due to the significant rise in turbine installation, condition monitoring activities have become more

challenging, which causes a substantial increase in operation and maintenance (O&M) costs. Therefore, many research activities are focusing on using advanced technologies to improve turbines' expensive components' life expectancy as well as minimising O&M costs. In addition to this, offshore wind farms (WFs) are generally situated in distant areas subject to harsh operating environmental conditions, which make offshore accessibility challenging and costly due to logistic and transport issues unlike onshore. A recent study [2] found that their O&M costs were estimated to account for 25–30% of the life cycle costs of an offshore WF. O&M visits make use of specialised vessels or helicopters for planned and unplanned repair activities. The O&M costs increase due to the higher incidence of failures that cause underperformance, high downtime and low availability. Furthermore, unscheduled maintenance occurs due to unexpected failures and affects the weather window; these need to be identified as quickly as possible to prevent critical damage and improve availability. All these factors together cause performance deterioration and significant loss in revenue that strongly affects the net economic value of the offshore WF [3]. Thus, WF developers and operators are continuously seeking cost-effective strategies to minimise O&M costs, improve performance, and maintain wind power availability while increasing at the same time, the return on their investment. According to recent statistics, it is estimated that the global business for WF O&M is expected to expand by up to \$ 27,400 USD million by 2025 [4], supported by data-driven technologies. This, therefore, outlines the importance of the O&M sector concerning the technical and commercial aspects for offshore WFs.

Condition monitoring (CM) is a process that has been widely used to monitor the operational status of machines in order to detect potential anomalies at an early stage to prevent catastrophic damage and to improve performance [5]. Three maintenance strategies are in use by the wind industry; these are reactive maintenance, corrective maintenance (run-to-failure) and preventive maintenance [5,6]. Preventive maintenance can be scheduled or condition-based, depending upon the problems. System physics-based techniques (e.g., oil debris analysis, vibration signal analysis) based on O&M activities consist of consolidation of run-to-failure and scheduled maintenance operations. Still, these techniques are costly as they cause significant downtime as well as a premature replacement of components. That is why O&M strategy is shifting from corrective and scheduled towards condition-based maintenance [6]. Predictive (condition-based) maintenance of the machine can be undertaken on a continuous basis without disrupting power generation, as well as being useful in determining the optimum point between corrective and scheduled maintenance. This improves maintenance activities and reduces unplanned downtime [7,8]. Within this framework, a condition-based maintenance approach is found to be cost-effective in improving reliability and reducing downtime [9,10], while minimising O&M costs [11,12].

As already stated, CM is an integral part of O&M that covers the essential activities associated with O&M at the different stages of WTs' operation. Turbine manufacturers and operators widely adopt supervisory control and data acquisition (SCADA) data-based CM as it improves performance and minimises O&M costs [13,14]. Qiu et al. [15] proposed a thermophysics-based approach (a synthesised thermal model) that uses SCADA data for WT drivetrain fault diagnosis by deriving relationships between various SCADA signals and changes in the thermophysics of WT operation. The results suggest that the SCADA-based thermophysics technique is useful in identifying non-linearity of the gearbox oil temperature rise with wind speed/output power, which can effectively suggest gearbox efficiency degradation that may be attributed to gear transmission problems such as gear teeth wear. Dao et al. [16] proposed a co-integration methodology based on SCADA data for improving CM and fault diagnosis that was found to be effective. Most of these technologies have concentrated on using SCADA signals rather than the SCADA alarms that are recorded in the SCADA system. Nevertheless, SCADA alarms are triggered and recorded when vital component signals exceed threshold limits. Hence, the inclusion of alarm signals can assist in identifying anomalies and improving WT performance. Recently, probability-based Venn diagrams and artificial neural networks [17–19], as well as classification methods [20], have been used for the analysis of SCADA alarm signals for WT condition/performance monitoring.

The WT power curve is widely used as a benchmark for use in purchase contracts and to correctly correlate the non-linear relationship between hub height wind speed (measured immediately before the rotor) and electrical power from the turbine [21]. Power curves can be used to optimise operational costs, enhance reliability and for condition/performance monitoring activities, and therefore are considered as critical performance indicators. However, power curves are adversely affected by changing environmental and topographical conditions, and thus may be site-specific [22,23]. Depending upon turbine rating and design, the commercial power curve shape deviates from the theoretical power curve [24,25]. Nevertheless, there has been extensive research on techniques for appraising power curves based on SCADA data over the last decades; these can be categorised into parametric and non-parametric methods as follows.

Parametric techniques include specified mathematical equations, perhaps from a family of functions, and with several parameters that are selected to provide the best fit to a particular WT power curve. Segmented linear models [26], polynomial regression [27,28], and models based on probabilistic distributions such as four- or five-parameter logistic distributions [29,30] have been used. Unlike parametric techniques, non-parametric approaches are adaptive. They can exhibit a high degree of flexibility because they do not enforce any pre-specified equation, and such methods can tightly fit the measured data subject to some specified smoothness of the fit criterion. Commonly used non-parametric techniques include ANN [31], SVM [32], GP [33] and Copula function [34]; they have proved to be useful in SCADA-based power curve modelling for improving WTs' forecasting and prediction as well as for O&M activities. In contrast to classical neural network techniques, SVM is useful in solving problems related to classification and prediction for non-linear issues. Vapnik proposed the SVM initially in 1992 [35], and it has been upgraded to provide better computational ability as well as higher prediction accuracy [36]. Santos et al. [37] proposed the SVM classification-based technique to identify failures related both to rotor blade imbalance and imbalance using simulated data points. The proposed algorithm compared different SVM kernels to neural networks with the conclusion that the linear kernel SVM outperforms other kernels and ANNs, in terms of validated accuracy, training and tuning times. Dahhani et al. [38] proposed an SVM-based control strategy for a WT, where SVM was used to detect optimal electromagnetic torque and blade pitch angle in response to wind speed changes. The results show that with just the knowledge of wind speed, SVM control could operate the overall wind power system optimally, which is validated by sliding mode control. Furthermore, SVM has also been applied in time-series wind speed forecasting [39], short-term wind power prediction [40] and CM [41,42] activities.

2. Scientific Novelty and Importance of This Research

The stochastic nature of wind and its complex interaction with the WT results in the variation of the power curve and significant uncertainty in its determination. This highlights the importance of uncertainty analysis associated with the power curve to assist turbine operators in the interpretation of performance validations. In forecasting, many articles such as [39,40,43], are presented to quantify the uncertainty for wind speed, wind power and smart grids purposes. De Brabanter et al. [44] proposed numerous techniques for calculating confidence intervals (CIs) for Least Squares Support Vector Regression. However, CI-based uncertainty assessment for SVM-based power curve models research is limited. Accurate estimation of uncertainty ensures early detection of anomalies caused by expected failures at an early stage and supports O&M decision management as highlighted by [45]. This research extends the work of [44] for SVM power curve construction in which two techniques, namely, simultaneous and pointwise, will be used as their CIs are close to the standard bootstrap method, compared to others. This paper fills this gap by proposing and comparing these CIs' techniques and suggesting which one is suitable for SVM-based power curve modelling.

Modern large WTs achieve peak values for C_p in the range of 45 to 50%. In addition, as seen from the above equations, WT mechanical power output depends on the power coefficient and the wind speed, while both blade pitch angle (β) and tip speed ratio (λ) have an impact on the power coefficient.

4. Data Description and Preparation

The SCADA data-based condition/performance monitoring is a cost-effective approach as it provides crucial information regarding the load history and operations of individual turbines. They provide an efficient tool for continuous CM of a turbine for early warning of failures and related performance issues. The SCADA data points collected from the Whitelee WF (located in Scotland, UK) are used for training and validating the proposed models, including data pre-processing, air density correction, and model evaluation. The Whitelee WF is located to the south of Glasgow, Scotland and amounts to 215 Siemens and Alstom onshore WTs with a total installed capacity of 539 MW. They record more than 100 different signals such as electrical (e.g., real and reactive power output and currents and voltages in the generator windings); weather-related signals (e.g., anemometer-measured hub height wind speed and direction, and ambient temperature); various temperature signals, such as main bearing and gearbox; pitch information (e.g., set and actual blade angles); numerous other signals. All this information is in the form of maximum, minimum, average and standard deviation values, with a 10 min average value. Measured wind speed is the most significant source of uncertainty and including more data points gives more certainty to the average value in the WT power curve. Type B uncertainties would be challenging to treat in a consistent manner without greater knowledge of the instrumentation used. Therefore, in this paper, we used the statistical spread for the SVM-based power curve using CIs. The data points used in this study cover the period from 00:00 on 1st September 2012 to 23:50 on 30th November 2012, accounting for 13,250 data samples in total. Recorded data can be imperfect due to sensor failures or malfunctions; these need to be removed or corrected as they affect the accuracy of any proposed models. Firstly, samples with missing values or negative power values are filtered out. Data points where maximum wind speed has reached more than 20 m/s are also filtered out because beyond this wind speed the turbine is stopped. In addition, data sampling during frequent start-up or stop in the low-wind-speed period may have a different variation. In short, criterion such as timestamp mismatch, negative power values, out of range values and turbine curtailment are used to filter out such misleading data similar to the one described in [13,23]. This reduces the number of data samples to 7918. Once the data are pre-processed, the next step is to carry out data partition into training and testing to make the model robust. Two problems need to be discussed when conducting data partition; with fewer training data, estimated values of the SVM model will have more significant variance while estimated values output statistics will have more substantial variance with less testing data. Thus, to guarantee variation within reasonable limits, a balance between training and testing data points is required. It should be noted that the inclusion of extensive training data improves data-driven models, such as SVM performance. In contrast, more testing data help in estimating errors accurately. Therefore, many studies [32–34] suggested that 70% training and 30% testing partition size are found to give the right combination of accuracy and precision. Hence, filtered SCADA data points are randomly divided into training and testing subsets in the ratio of 70:30, i.e., 5542 for training and 2376 for evaluation, as illustrated in Table 1. The SVM model has never seen the testing data points, so the resulting outcome will be an excellent guide to analyse its impact on power curve accuracy and uncertainty when the model is applied to unseen data and this is discussed in the sections below.

Table 1. Supervisory control and data acquisition (SCADA) data description.

Start Timestamp	End Timestamp	Measured Data	Filtered Data	Training Data	Validation Data
101/09/2012 00:00 A.M.	30/11/2012 23:50 P.M.	13,250	7918	5542	2376

Air Density Correction for Improving Power Curve Accuracy

The International Electrotechnical Commission (IEC) standard 61400-12-1 [46–48] proposes a technique for variable pitch regulated WTs to compensate for the air density effect on a power curve, where atmospheric pressure and ambient temperature are acknowledged as relevant parameters for air density correction. Therefore, the following equations (defined in the IEC standard) are used in this study for air density correction purposes:

$$\rho = 1.225 \left[\frac{288.15}{T} \right] \left[\frac{B}{1013.3} \right] \tag{3}$$

and

$$V_C = V_M \left[\frac{\rho}{1.225} \right]^{\frac{1}{3}} \tag{4}$$

where V_M and V_C are the measured and corrected wind speed (m/sec), respectively. Ambient temperature (T) and atmospheric pressure (B) average 10 min SCADA data points are used in Equation (3) to calculate the corrected air density and then the corrected wind speed (V_C) using Equation (4). The corrected wind speed is then plotted against electrical power output to provide the corrected power curve. IEC air density correction is widely used in data-driven models to construct the power curve.

Figure 2 shows the measured power curve of a turbine before filtering and air density correction. Adopting air density correction (as described in Section 3) and filtering, as outlined above, results are illustrated in Figure 3. The changes due to air density correction is not significant as data used in this study came from wind turbines from normal region where temperature is not significant [47]. It is worth noting that WT power curves are affected by both environmental and operational conditions and incorporating those could be useful for improving data-driven power curve accuracies [47,48].

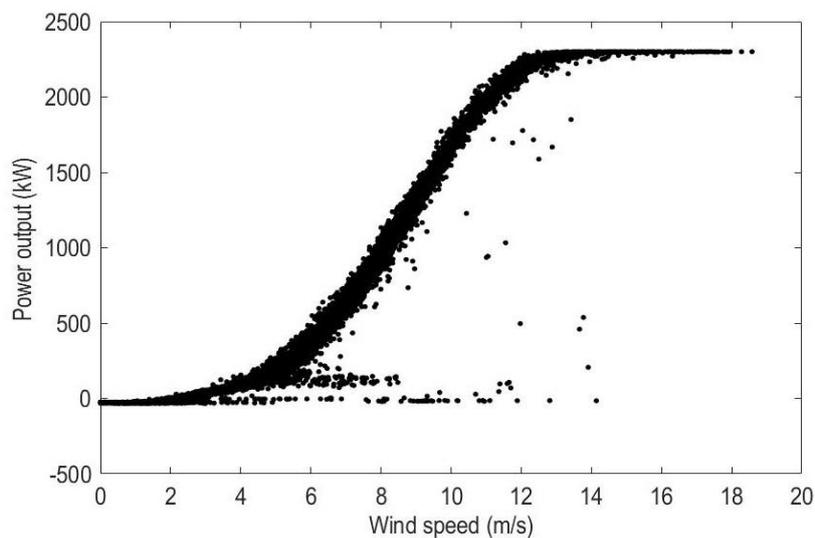
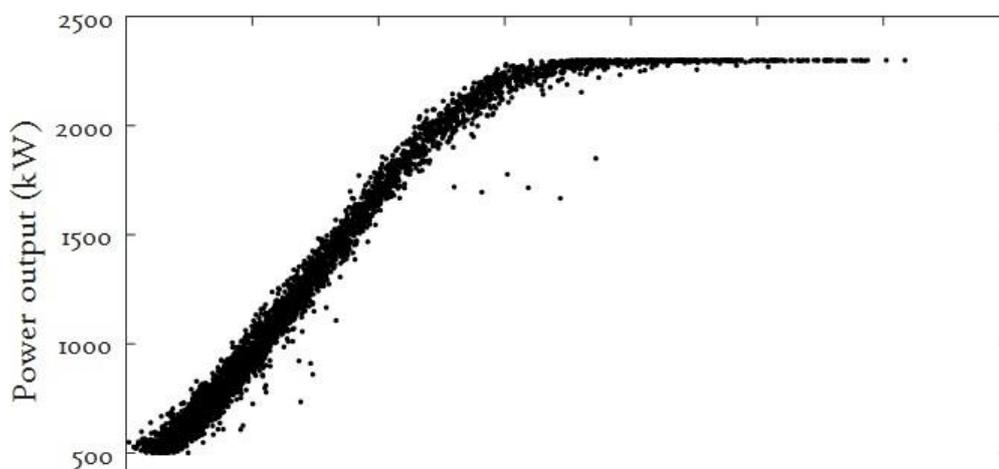


Figure 2. Raw power curve data.



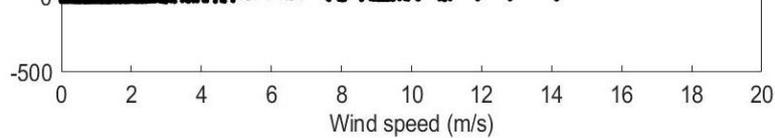


Figure 2. Raw power curve data.

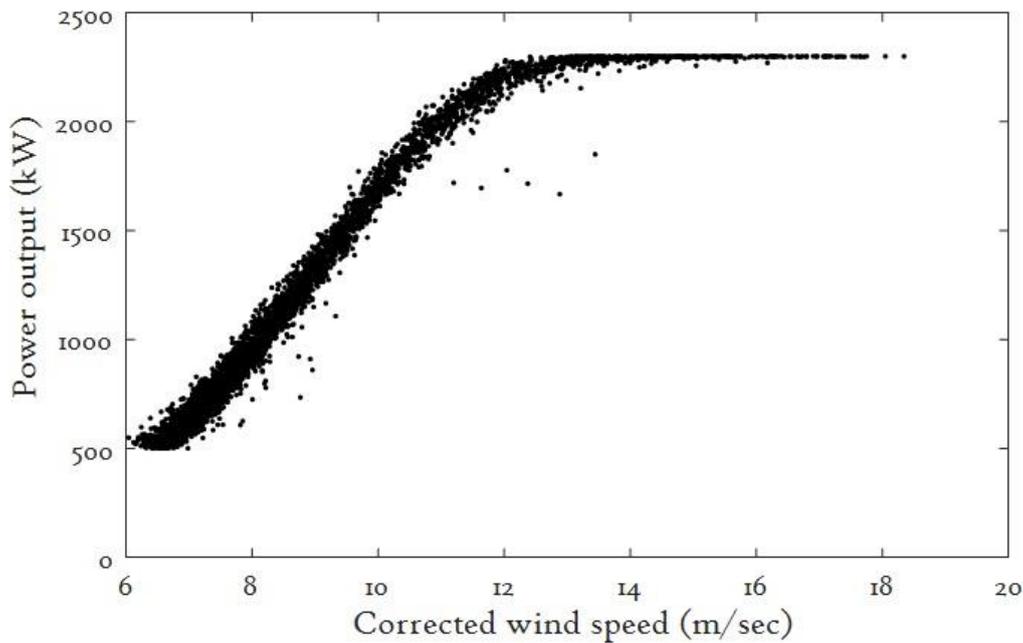


Figure 3. Air density corrected filtered power curve.

5. Methodologies

5.1. SVM Models—Theoretical Descriptions

The SVM is a non-linear, data-driven technique, gaining popularity due to superior performance as compared to traditional ERM as used by conventional neural networks [49]. With the use of SRM, the upper bound on the anticipated risk minimises the effect of a reduction in training data error and gives SVM greater capacity to generalise function as compared to neural networks. Initially, SVMs were constructed to provide objective/optimal classification called SVC but more recently have been applied to regression and termed as SVR. In this section, the theoretical description of SVM regression models [36,50] for WT power curves is described as follows.

Let us consider N training vectors $x_i \in \mathfrak{R}$ defined by a set of definite variables $x_i = \{x_{i1}, x_{i2}, x_{i3}, \dots, x_{ip}\}$ and by the class response $y_i \in \mathfrak{R}$. To use non-linear functions for regression of the data x , a non-linear map $\emptyset: x \rightarrow \emptyset(x)$ into a high dimensional space is suggested to allow linear regression in that space ($f(x) = \langle w, \emptyset(x) \rangle + b$). Dot products of $\langle x_i, x \rangle$ found in calculating linear SVR are in the non-linear case replaced by the dot products $\langle \emptyset(x_i), \emptyset(x) \rangle$. This is the symmetric function in x_i and x that should satisfy Mercer’s condition and is called a kernel: $(x_i, x) = \langle \emptyset(x_i), \emptyset(x) \rangle$. The SVR technique is formulated as a minimisation of the following function:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \tag{5}$$

subject to:

$$\|f(x_i) - y_i\| \leq \varepsilon, \quad i = 1, 2, \dots, N \tag{6}$$

where

$$f(x_i) = \langle w, \emptyset(x_i) \rangle + b \tag{7}$$

Here, w is a space weight coefficient vector, and b is a bias.

The SVM is a supervised machine learning technique whose sole goal is to design a hyperplane that classifies all training vectors into two classes. The most appropriate hyperplane is one that leaves the maximum margin between both classes. Therefore, minimising the term w maximises the separability. Minimisation of w is a non-linear, optimisation task, which can be solved by the

Karush–Kuhn–Tucker (KKT) approach [50], using Lagrange function/multipliers to find a ξ -insensitive loss function, by the following function:

$$f(x) = \sum_{i=1}^N (\alpha_i - \alpha_i^*) K(x_i, x) + b \tag{8}$$

with the property:

$$w = \sum_{i=1}^N (\alpha_i - \alpha_i^*) \varphi(x_i) \tag{9}$$

where b is the bias term (a scalar), and α_i and $\alpha_i^* \geq 0$ are the Lagrange multipliers. The sample point that appears with non-zero coefficients α_i is called the support vector.

A slack variable C is introduced into (5) to generate an optimal SVR model by relaxing the margin constraints and then neglecting a controlled part of the data. This allows the optimisation problem to be reformulated as Equation (10) below,

$$\min_{w, b, \xi^-, \xi^+} \frac{1}{2} \|w\|^2 + C \sum_{k=1}^N (\xi_k^- + \xi_k^+) \tag{10}$$

subject to:

$$y_i - w^T \varphi(x_i) - b \leq \varepsilon + \xi_k^-, \quad i = 1, 2, \dots, N$$

and

$$-y_i + w^T \varphi(x_i) + b \leq \varepsilon + \xi_k^+, \quad i = 1, 2, \dots, N$$

where $\xi_k^-, \xi_k^+ \geq 0$ are slack variables that cause a penalty term, which is weighted by C and are used to measure the deviations of the samples outside the ξ -insensitive zone. The addition of a slack variable lies in the following range: $0 \leq \alpha_i, \alpha_i^* \leq C$. Weight of misclassifications increases with the increase in C values, which leads to a higher cost of the misclassified data points that cause strict separation of data. This factor C is called a box constraint because it is in the formulation of the dual optimisation problem where the Lagrange multipliers are bounded within the range $[0, C]$. Minimising the first term of Equation (10) requires that the function fitting through data be as flat as possible; minimising the second term penalises deviations more significant than ξ , which is tuned by C .

The Gaussian kernel, radial basis function (RBF) with the kernel scale σ , is used in this study for data mapping as it facilitates computations in higher-dimensional space in a better way. Mathematically the RBF is defined by:

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) \tag{11}$$

5.2. Uncertainty Estimation–Theoretical Descriptions

These CIs are often estimated for uncertainty analysis and for providing turbine operators with confidence in determining how well a particular model describes the actual underlying process by taking into account the estimator’s statistical properties. As discussed earlier, CIs are also vital for identifying anomalies associated with the malfunction of the turbine. Additionally, data points that are not within a specified CI range can be considered anomalous and are potentially caused by damage to the turbine [42]. De Brabanter et al.’s [44] work is extended here to construct two CI techniques, namely, pointwise and simultaneous methods, for measuring uncertainty associated with the SVM power curve [44,48]; these are briefly described as follows.

The difference between the average estimate of the model and the average measured value determines the “bias”. At the same time “variance” is used to express the variability of the model estimation for given data points and indicate the spread of data. Both these are measures of model prediction accuracy. For example, high bias resulted in a high error on training as well as testing data and resulted in bias estimation problems. Therefore, the bias estimation problem needs to be taken care of by CIs so that the interval is correctly centred as well as providing proper coverage [51]. CIs are

based on the central limit theorem for linear smoothers combined with bias correction and variance estimation. The following equations are used in this study to determine bias and variance [44]:

$$\hat{b}[\hat{m}(x)|X = x] = L(x)^T \hat{m} - \hat{m}(x) \tag{12}$$

$$Var[\hat{m}(x)|X = x] = L(x)^T \hat{\Sigma}^2 L(x) \tag{13}$$

with

$$\hat{\Sigma}^2 = diag(\hat{\sigma}^2(X_1), \dots, \hat{\sigma}^2(X_n))$$

where $L(x)$ is the smoother vector evaluated at a point x and denoted $\hat{m} = (\hat{m}(X_1), \dots, \hat{m}(X_n))^T$.

The residuals are calculated from the SVM-based power curve to determine conditional bias and variance. These will later be used for estimating CIs for SVM power curve uncertainty analysis.

5.2.1. Pointwise CIs for modelled SVM power curve

Pointwise CIs for the SVM power curve model are calculated by the following equation [44]:

$$\hat{m}(x) \pm z_{1-\frac{\alpha}{2}} \sqrt{Var[\hat{m}(x)|X = x]} \tag{14}$$

where $z_{1-\alpha/2}$ is the $(1 - \alpha/2)$ th quantile of the standard Gaussian distribution. In Equation (14), plus and minus signs signify upper and lower estimated CI values, respectively. However, this equation excludes bias, and so provides an inaccurate picture of uncertainty associated with the SVM model. Hence, unknown bias is estimated by Equation (12) which is incorporated in Equation (14) to reflect bias. Hence, the modified formula for pointwise CIs that includes a bias-corrected approximation 100(1 - α)% is expressed by the following equation:

$$\hat{m}_c(x) \pm z_{1-\alpha/2} \sqrt{Var[\hat{m}(x)|X = x]} \tag{15}$$

where $\hat{m}_c(x) = \hat{m}(x) - \hat{b}[\hat{m}(x)|X = x]$ and α is called the significance level and kept at 0.05, corresponding to 95% probability. The corrected bias values are further normalised to determine the $(1 - \alpha/2)$ th quantile of the standard Gaussian distribution.

5.2.2. Simultaneous CIs for Modelled SVM Power Curve

Simultaneous CIs are defined as intervals that constitute specific intervals for the independent components of the parameter and are advantageous (in terms of computation and mathematical complexity) for multiple comparisons that include the combination of several single CIs, in contrast to multiple CIs. In this study, simultaneous CIs are constructed for the SVM-based power curve and then compared with pointwise CIs to determine which approach is most robust. Several studies have published novel methods to calculate simultaneous CIs such as [51,52]. However, the Bonferroni and Sidak corrections techniques were found to be accessible as they are mathematically easy to implement and produce acceptably accurate results, see, for example [53]. Hence, Equation (15) is modified based on the Sidak correction to determine simultaneous CIs, and is expressed by the following equation [50]:

$$\hat{m}_c(x) \pm z_{1-\beta/2} \sqrt{Var[\hat{m}(x)|X = x]} \tag{16}$$

subject to: $\hat{m}_c(x) = \hat{m}(x) - \hat{b}[\hat{m}(x)|X = x]$. β is assumed as the significance threshold for each test and determined by $\beta = 1 - (1 - \alpha)^{1/n}$ where n is the number of test samples. Equation (16) is further modified to include the approximate 100(1 - α)% and is shown below,

$$\hat{m}_c(x) \pm v_{1-\alpha} \sqrt{Var[\hat{m}(x)|X = x]} \tag{17}$$

where $v_{1-\alpha} = \sqrt{2 \log\left(\frac{k_0}{\alpha\pi}\right)}$ and subject to:

$$k_0 = \int_0^\chi \frac{\sqrt{\|L(x)\|^2 \|L'(x)\|^2 - (L(x)^T L'(x))^2}}{\|L(x)\|^2} dx$$

where χ denotes the set of x values of interest and $L'(x) = \left(\frac{d}{dx}\right)L(x)$ in which elementwise differentiation is applied. The k_0 is strongly related to degrees of freedom of the fit and approximates to the following relationship [44]: $k_0 \approx \frac{\pi}{2}(\text{trace}(L) - 1)$. All these values are calculated and used with Equation (17) to determine simultaneous CIs for the SVM model.

6. Results and Discussions

6.1. SVM-Based Power Curve Model

The pre-processed and air density corrected training SCADA data points (described in Section 6) were used to train the SVR model based on the above-described methodology and then testing datasets were used to validate the effectiveness of the model. The training and testing datasets of Table 1 are used to train and validate the SVM power curve model where wind speed used as input to estimate the power which is then plotted together as shown in Figure 4. It is worth noting that the Matlab “fitsvm” function with the “OptimizeHyperparameters” with “automatic” options is used to calculate the “optimal value” for the BoxConstraint [54], while the K-fold CV approach, as per methodology described in [55], is used to calculate the value of kernel width/scale. The SVM-based power curve model intrinsically represents fitting errors. The SVR-based power curve model is found to be continuous, smooth and accurately estimates the measured power curve below the rated power. This is further confirmed by plotting estimated and measured power values in a time-series, as shown in Figure 5, and is plotted with limited datasets for better visualisation of the results. However, SVM model accuracy depends on the data quantity and quality, as well as the appropriate method of fitting used. Hence, the performance of the SVM model observed over 16 months suffers because of the lack of quality data points.

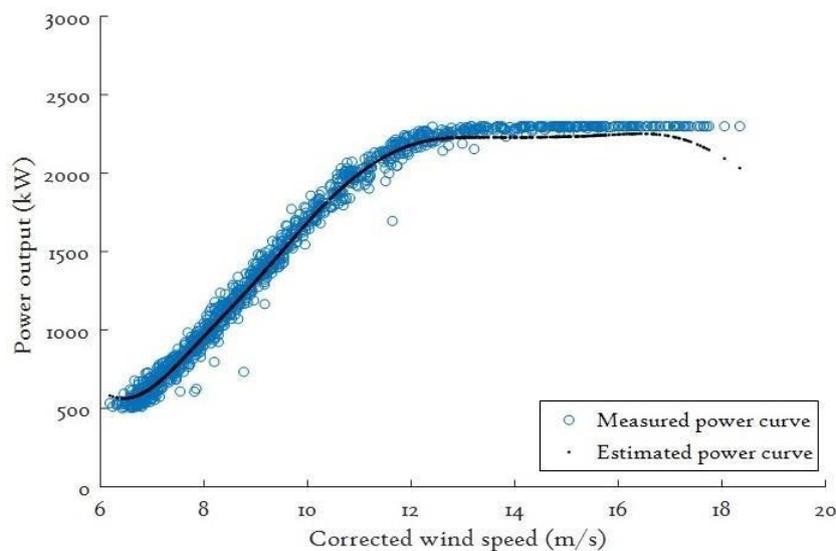
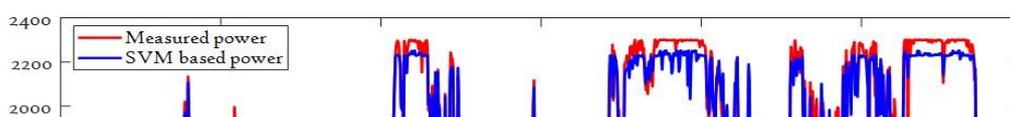


Figure 4. Estimated power curve and measured data.

The differences between the observed value of the dependent variable and the estimated value are called residuals, and they are vital in determining the deviation between measured data and the regression model. The frequency distribution of the calculated residuals of the SVR power curve is plotted in Figure 6, together with a fitted Gaussian distribution and the results suggest that the distribution of SVR residuals closely follows a Gaussian distribution.



The differences between the observed value of the dependent variable and the estimated value are called residuals, and they are vital in determining the deviation between measured data and the regression model. The frequency distribution of the calculated residuals of the SVR power curve is plotted in Figure 6, together with a fitted Gaussian distribution and the results suggest that the distribution of SVR residuals closely follows a Gaussian distribution.

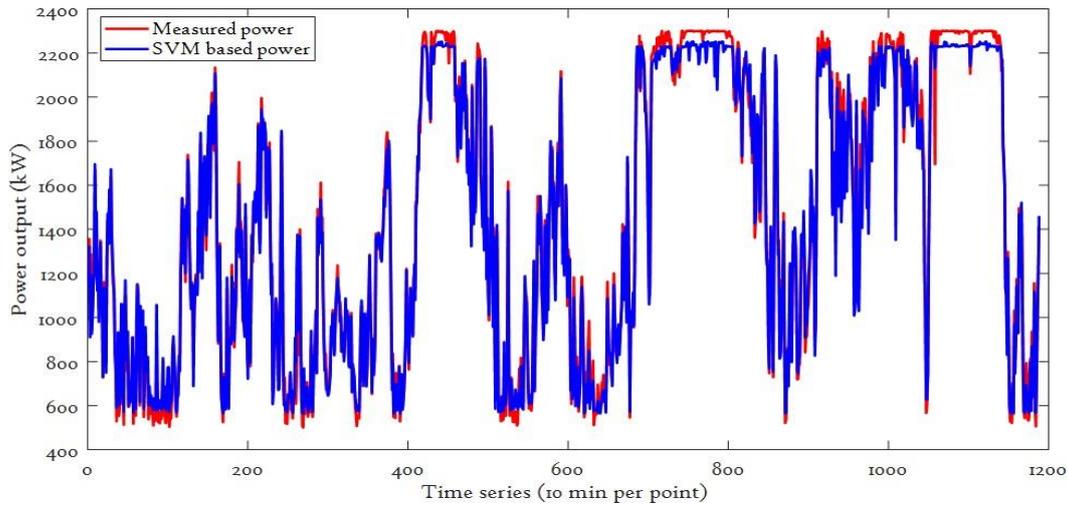


Figure 5. Comparative studies of SVM power curve model in terms of time-series.

The differences between the observed value of the dependent variable and the estimated value are called residuals, and they are vital in determining the deviation between measured data and the regression model. The frequency distribution of the calculated residuals of the SVR power curve is plotted in Figure 6, together with a fitted Gaussian distribution and the results suggest that the distribution of SVR residuals closely follows a Gaussian distribution.

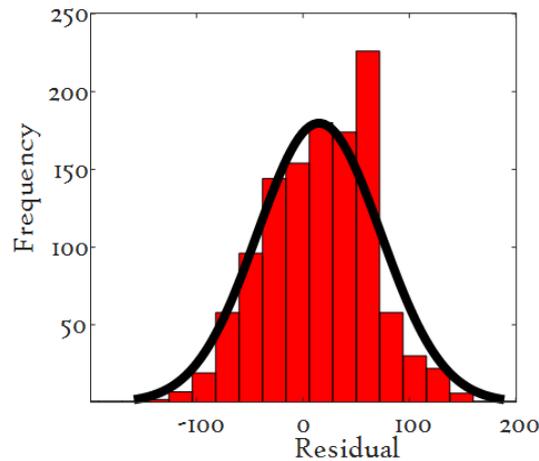


Figure 6. Histogram fitting of SVM-based power curve model.

Furthermore, statistical performance indices, namely, root-mean-squared error (RMSE), mean absolute error (MAE), and coefficient of determination (R^2) are used to measure the performance of the SVR fitting. The calculated values of RMSE (59.93), MAE (46.18), and R^2 (0.99) further confirm the accuracy of the SVR power curve model.

6.2. SVM Power Curve Uncertainty Analysis

As stated in previous sections, CIs provide information on the uncertainty surrounding an estimation but are themselves model-based estimates that reflect the standard deviation of the model, for example, see Equations (14)–(17), and therefore can be valuable for early fault detection as demonstrated by [42]. Therefore, using the methodologies described above, CIs are calculated for analysing SVM-based power curve uncertainties. After that, at the later stage, a comparative study among these techniques is carried out to suggest which technique is robust and accurately estimates the uncertainty associated with the SVM power curve.

Using Equation (15), pointwise CIs calculated for the SVM power curve and plotted together with the estimated mean and measured values, as shown in Figure 7, suggest that the estimated and measured power curves are mostly within the region defined by the pointwise CIs. To understand this better, measured and estimated values are plotted in a time-series plot for selected power range and illustrated in Figure 8, which depicts that CIs upper and lower bounds have a tight width within

Using Equation (15), pointwise CIs calculated for the SVM power curve and plotted together with the estimated mean and measured values, as shown in Figure 7, suggest that the estimated and measured power curves are mostly within the region defined by the pointwise CIs. To understand this better, measured and estimated values are plotted in a time-series plot for selected power range and illustrated in Figure 8, which depicts that CIs upper and lower bounds have a tight width within predicted power values. Similarly, using Equation (17), simultaneous CIs calculated for the SVM power curve, as shown in Figure 9, signify that the estimated and measured power curves follow the expected variance of the measured data. It should be noted that in Figures 8 and 10, selected time-series data have been used to avoid complex figures and to explain the results in a better way. However, they have a large width across all the wind speed range, which is clearly seen by plotting time-series between measured and estimated values. In Figure 9, the simultaneous CIs' bandwidth starts to get wider near to a wind speed value of 13 m/sec, and the time-series values of power between 700–800 in Figure 10 demonstrate this further.

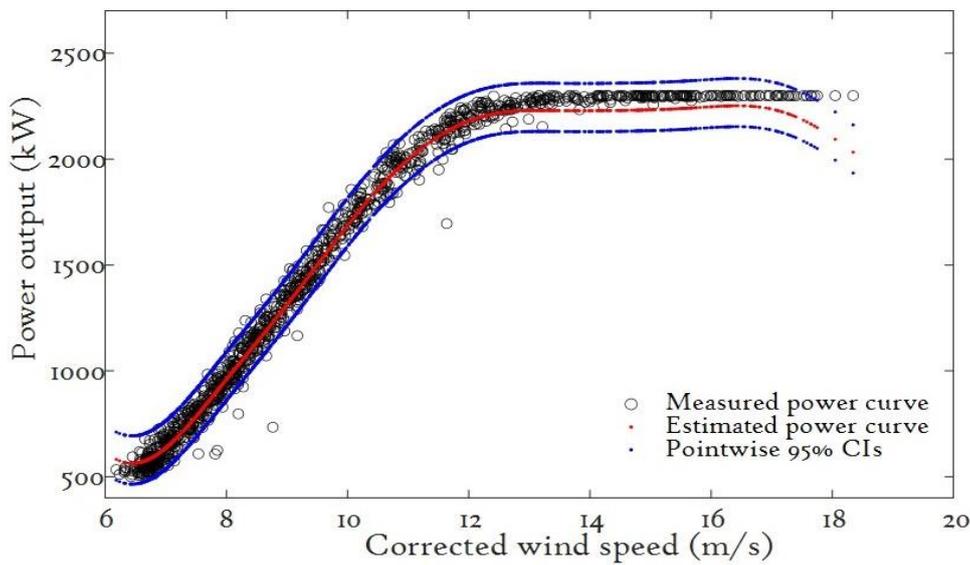


Figure 7. Pointwise CIs for SVM-based power curve.

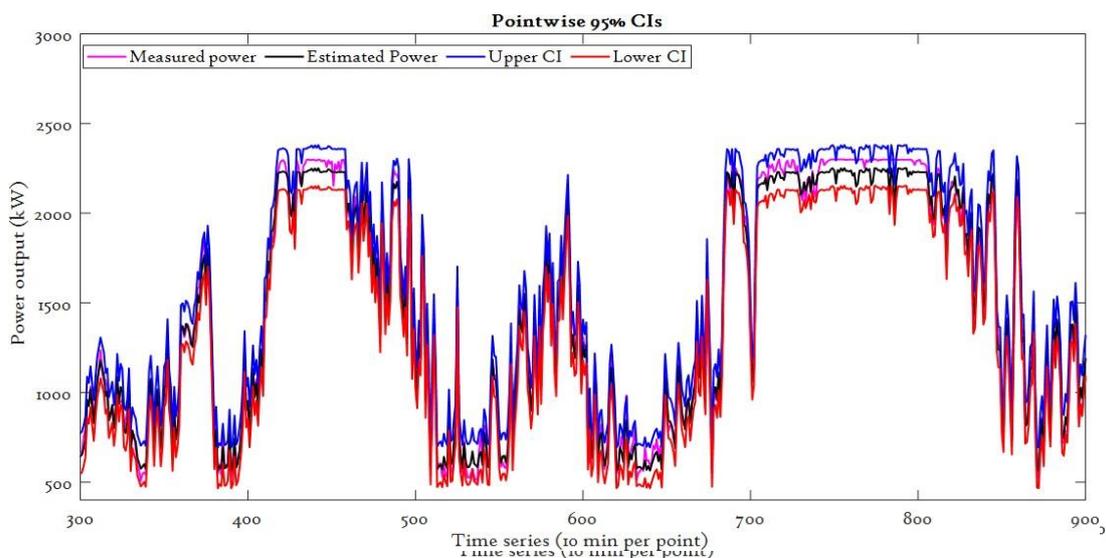
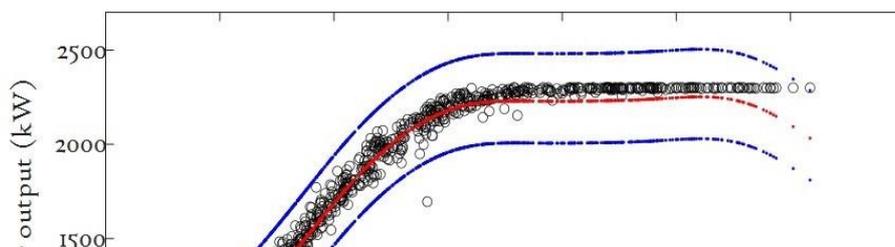


Figure 8. Analysis of pointwise CIs for estimated power values in time-series.



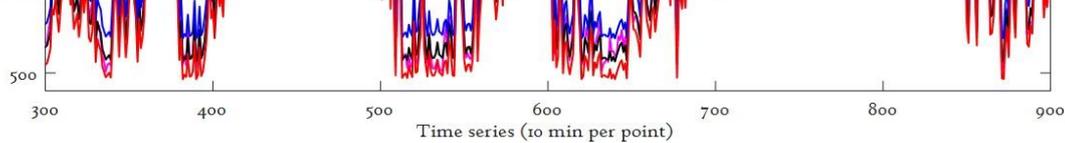


Figure 8. Analysis of pointwise CIs for estimated power values in time-series.

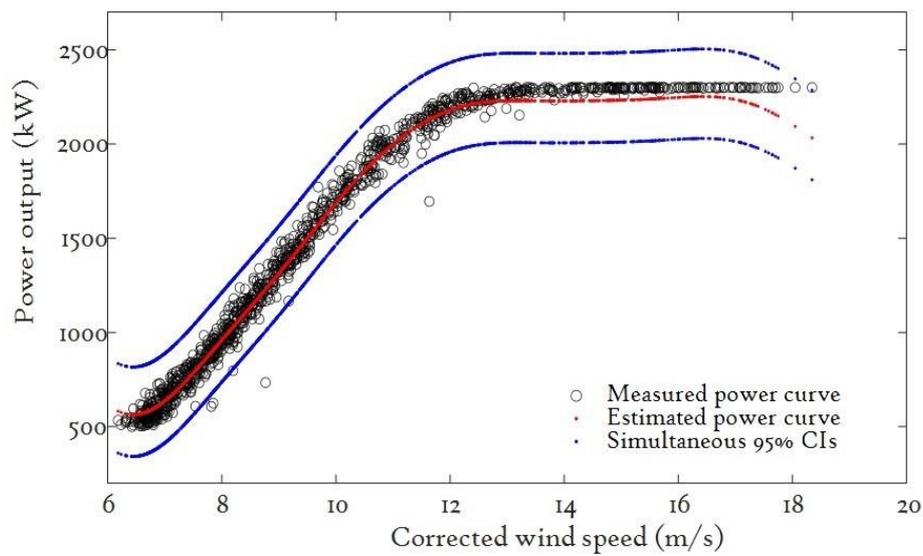


Figure 9. Simultaneous CIs for SVM-based power curve.

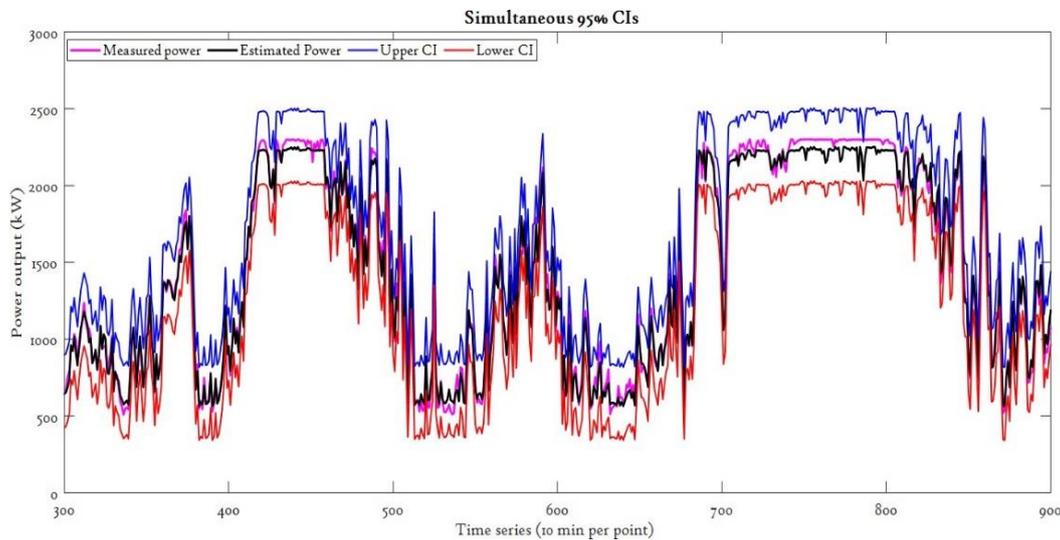


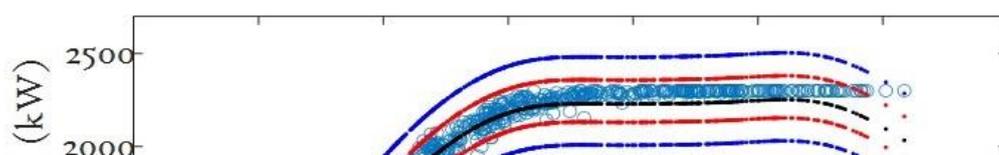
Figure 10. Analysis of simultaneous CIs for estimated power values in time-series.

6.3. Comparative Studies of the Proposed Methodologies

The confidence interval of the SVM power curve, being smaller for the critical wind speed range (between cut-in and cut-out), may have better ability to reject the anomalous or faulty data, and thus be helpful in optimising WTs, and have more remarkable ability to detect anomalies in the early stages; this section addresses this as follows.

The developed SVM power curve, together with estimated pointwise and simultaneous CIs' results are, along with the measured power curve, shown in Figure 11. They suggest that pointwise CIs are relatively smaller across the entire range of wind speed and therefore have a superior capability to reject invalid or faulty data as compared to simultaneous CIs.

This difference can be further seen in the limited time-series plot (Figure 12) where dashed circles highlight significant smaller pointwise CIs and thus have reduced uncertainty for the SVM power curve, as compared to simultaneous CIs. It should be noted that the SVM model's accuracy weakens in dealing with extensive datasets due to cubic inversion issues and, therefore, dealing with a large data size can be challenging and consequently affects the proposed uncertainty model's accuracy. Therefore, finding appropriate data management is vital for the effective use of SVM models.



highlight significant smaller pointwise CIs and thus have reduced uncertainty for the SVM power curve, as compared to simultaneous CIs. It should be noted that the SVM model's accuracy weakens in dealing with extensive datasets due to cubic inversion issues and, therefore, dealing with a large data size can be challenging and consequently affects the proposed uncertainty model's accuracy. Therefore, finding appropriate data management is vital for the effective use of SVM models. 14 of 18

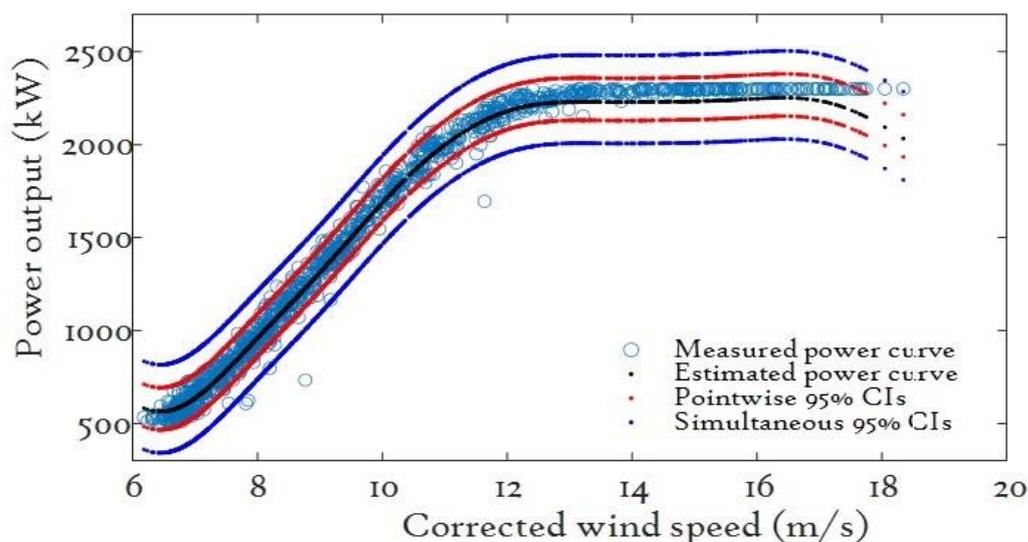


Figure 11. SVM-based power curve comparative investigation.
Figure 11. SVM-based power curve comparative investigation.

This difference can be further seen in the limited time-series plot (Figure 12) where dashed circles highlight significant smaller pointwise CIs and thus have reduced uncertainty for the SVM power curve, as compared to simultaneous CIs. It should be noted that the SVM model's accuracy weakens in dealing with extensive datasets due to cubic inversion issues and, therefore, dealing with a large data size can be challenging and consequently affects the proposed uncertainty model's accuracy. Therefore, finding appropriate data management is vital for the effective use of SVM models.

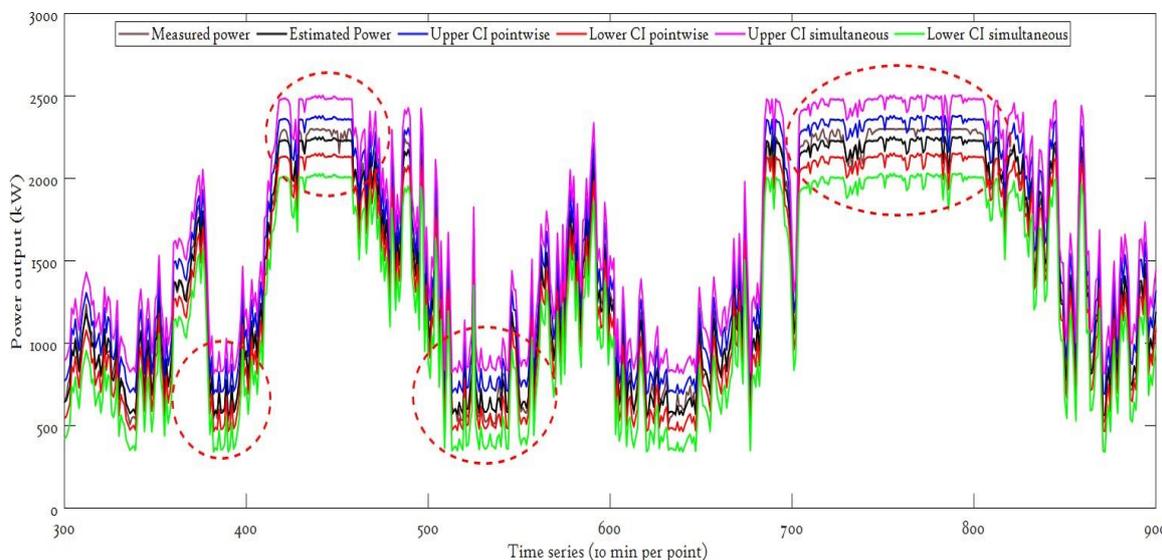


Figure 12. Comparative studies on SVM models in terms of time-series.

7.7. Conclusions

Data driven power curves are widely used by both turbine operators and service for performance monitoring and O&M mainly to formulations related to WTs. However, its accurate quantification related to the power curve remains a challenging issue. In this paper, two approaches are proposed to estimate the SVM based power SVM based uncertainty model with their merits and demerits. SCADA datasets to SCADA data health operational pitch regulated WTs are used to train WTs and test the SVM model and the SVM uncertainty. Both pointwise CIs and simultaneous CIs are formulated to be used in estimating the uncertainty of the SVM power curve. The SVM power curve and the SVM power curve in Figures 7 and 9. However, by comparing simultaneously with pointwise CIs, it has been found that the latter is more accurate over the entire section of the SVM power curve, as illustrated in Figures 11 and 12. This is because pointwise CIs show a relatively narrow width across the entire wind speed range and therefore have better ability to detect anomalies at an early stage and improve WTs' maintenance decision process and other relevant activities, as compared to simultaneous CIs. Therefore, future research will extend this work by developing improved SVM-based failure detection for WTs' CM

by comparing simultaneously with pointwise CIs, it has been found that the latter is more accurate over the entire section of the SVM power curve, as illustrated in Figures 11 and 12. This is because pointwise CIs show a relatively narrow width across the entire wind speed range and therefore have better ability to detect anomalies at an early stage and improve WTs' maintenance decision process and other relevant activities, as compared to simultaneous CIs. Therefore, future research will extend this work by developing improved SVM-based failure detection for WTs' CM, where pointwise CIs (obtained from this study) will be used to identify the early signs of failure. In addition to this, studying the impact of environmental and operational conditions on the SVM power curve models' accuracy and uncertainty is also kept for future investigation.

Author Contributions: Conceptualisation, methodology and software, R.P.; validation, formal analysis and investigation, R.P. and A.K.; writing—original draft preparation, R.P.; writing—review and editing, A.K. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Acknowledgments: This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement no. 74 5625 (ROME0) ("Romeo Project" 2018). The dissemination of results herein reflects only the authors' views, and the European Commission is not responsible for any use that may be made of the information the paper contains.

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

Nomenclature

ANNs	Artificial neural networks
CIs	Confidence intervals
ERM	Empirical risk minimisation
GP	Gaussian process
GWEC	Global Wind Energy Council
SVM	Support vector machine
SRM	Structural risk minimisation
SVC	Support vector classification
SVR	Support vector regression
K	The general covariance matrix
KKT	Karush–Kuhn–Tucker
SCADA	Supervisory control and data acquisition
σ	Kernel width/scale
ξ	Insensitive zone
C	Box constraint
α	Significance level
B	Bias
MAE	Mean absolute error
R^2	Coefficient of determination
RMSE	Root-mean-squared error
RBF	Radial basis function
W	Space weight coefficient vector
WTs	Wind turbines
ξ_k^-, ξ_k^+	Slack variables

References

1. GWEC Report Entitled 'GWEC Forecasts 817 GW of Wind Power in 2021'. Available online: <https://gwec.net/global-wind-report-2019/> (accessed on 20 October 2020).
2. Ioannou, A.; Angus, A.; Brennan, F. A life-cycle techno-economic model of offshore wind energy for different entry and exit instances. *Appl. Energy* **2018**, *221*, 406–424. [CrossRef]

3. Honrubia-Escribano, A.; Martín-Martínez, S.; Honrubia-Escribano, A.; Gómez-Lázaro, E. Wind turbine reliability: A comprehensive review towards effective condition monitoring development. *Appl. Energy* **2018**, *228*, 1569–1583.
4. Agarwal, A. Wind turbine operations and maintenance market—Global market size, trends, and key country analysis to 2025. *Technol. Rep. Glob. Data* **2017**. Available online: <https://www.pes.eu.com/wind/global-wind-operations-and-maintenance-market-set-to-hit-27-4-billion-by-2025-says-globaldata/#:~:text=The%20global%20wind%20operations%20and,research%20and%20consulting%20firm%20GlobalData> (accessed on 20 October 2020).
5. Kolios, A.; Walgern, J.; Koukoura, S.; Pandit, R.; Chiachio-Ruano, J. Open O&M: Robust O&M open access tool for improving operation and maintenance of offshore wind turbines. In Proceedings of the 29th European Safety and Reliability Conference (ESREL), Hannover, Germany, 22–26 September 2019; pp. 452–459.
6. Scheu, M.N.; Tremps, L.; Smolka, U.; Kolios, A.; Brennan, F. A systematic Failure Mode Effects and Criticality Analysis for offshore wind turbine systems towards integrated condition based maintenance strategies. *Ocean Eng.* **2019**, *176*, 118–133. [[CrossRef](#)]
7. Qiao, W.; Zhang, P.; Chow, M.-Y. Condition monitoring, diagnosis, prognosis, and health management for wind energy conversion systems. *IEEE Trans. Ind. Electron.* **2015**, *62*, 6533–6535. [[CrossRef](#)]
8. Tian, Z.; Jin, T.; Wu, B.; Ding, F. Condition based maintenance optimisation for wind power generation systems under continuous monitoring. *Renew. Energy* **2011**, *36*, 1502–1509. [[CrossRef](#)]
9. Bussell, G.J.W.; Zaaijer, M.B. Reliability, availability and maintenance aspects of large-scale offshore wind farms, a concepts study. In Proceedings of the MAREC 2001: 2-day International Conference on Marine Renewable Energies, London, UK, 2001; Volume 113, p. 226.
10. Lu, B.; Li, Y.; Wu, X.; Yang, Z. A review of recent advances in wind turbine condition monitoring and fault diagnosis. In Proceedings of the IEEE Power Electronics and Machines in Wind Applications (PEMWA), Lincoln, Nebraska, 24–26 June 2009; pp. 1–7.
11. Qian, P.; Ma, X.; Zhang, D.; Wang, J. Data-Driven Condition Monitoring Approaches to Improving Power Output of Wind Turbines. *IEEE Trans. Ind. Electron.* **2018**, *66*, 6012–6020. [[CrossRef](#)]
12. Moeini, R.; Entezami, M.; Ratkovac, M.; Tricoli, P.; Hemida, H.; Hoeffler, R.; Baniotopoulos, C. Perspectives on condition monitoring techniques of wind turbines. *Wind Eng.* **2019**, *43*, 539–555. [[CrossRef](#)]
13. Bangalore, P.; Patriksson, M. Analysis of SCADA data for early fault detection, with application to the maintenance management of wind turbines. *Renew. Energy* **2018**, *115*, 521–532. [[CrossRef](#)]
14. Herp, J.; Pedersen, N.L.; Nadimi, E.S. A Novel Probabilistic Long-Term Fault Prediction Framework beyond SCADA. *J. Phys. Conf. Ser.* **2019**, *1222*, 012043. [[CrossRef](#)]
15. Qiu, Y.; Feng, Y.; Sun, J.; Zhang, W.; Infield, D. Applying thermophysics for wind turbine drivetrain fault diagnosis using SCADA data. *IET Renew. Power Gener.* **2016**, *10*, 661–668. [[CrossRef](#)]
16. Dao, P.B.; Staszewski, W.J.; Barszcz, T.; Uhl, T. Condition monitoring and fault detection in wind turbines based on co-integration analysis of SCADA data. *Renew. Energy* **2018**, *16*, 107–122. [[CrossRef](#)]
17. Qiu, Y.; Feng, Y.; Infield, D. Fault diagnosis of wind turbine with SCADA alarms based multidimensional information processing method. *Renew. Energy* **2020**, *145*, 1923–1931. [[CrossRef](#)]
18. Qiu, Y.; Feng, Y.; Tavner, P.; Richardson, P.; Erdős, F.G.; Chen, B. Wind turbine SCADA alarm analysis for improving reliability. *Wind. Energy* **2011**, *15*, 951–966. [[CrossRef](#)]
19. Chen, B.; Qiu, Y.; Feng, Y.; Tavner, P.; Song, W. Wind turbine SCADA alarm pattern recognition. In Proceedings of the IET Conference on Renewable Power Generation (RPG 2011), Edinburgh, UK, 6–8 September 2011.
20. Leahy, K.; Gallagher, C.; O'Donovan, P.; Bruton, K.; O'Sullivan, D.T.J. A robust prescriptive framework and performance metric for diagnosing and predicting wind turbine faults based on SCADA and alarms data with case study. *Energies* **2018**, *11*, 1738. [[CrossRef](#)]
21. Burton, T.; Sharpe, D.; Jenkins, N.; Bossanyi, E. *Wind Energy Handbook*; Wiley-Blackwell: Hoboken, NJ, USA, 2011; ISBN 0471489972. [[CrossRef](#)]
22. Dupré, A.; Drobinski, P.; Badosa, J.; Briard, C.; Plougonven, R. Air Density Induced Error on Wind Energy Estimation. *Ann. Geophys. Discuss.* **2019**. [[CrossRef](#)]
23. Koukoura, S.; Carroll, J.; McDonald, A. An insight into wind turbine planet bearing fault prediction using SCADA data. In Proceedings of the European Conference of the PHM Society, Utrecht, The Netherlands, 30 January 2018; Volume 4.

24. Thapar, V.; Agnihotri, G.; Sethi, V.K. Critical analysis of methods for mathematical modelling of wind turbines. *Renew. Energy* **2011**, *36*, 3166–3177. [[CrossRef](#)]
25. Dongre, B.; Pateriya, R.K. Adaptive filter-based power curve modeling to estimate wind turbine power output. *Wind. Eng.* **2019**. [[CrossRef](#)]
26. Khalfallah, M.; Koliub, A. Suggestions for improving wind turbines power curves. *Desalination* **2007**, *209*, 221–229. [[CrossRef](#)]
27. Marčiukaitis, M.; Žutautaitė, I.; Martišauskas, L.; Jokšas, B.; Gecevičius, G.; Sfetsos, A. Non-linear regression model for wind turbine power curve. *Renew. Energy* **2017**, *113*, 732–741. [[CrossRef](#)]
28. Raj, M.M.; Alexander, M.; Lydia, M. Modeling of wind turbine power curve. In Proceedings of the ISGT 2011, Kollam, India, 1–3 December 2011; pp. 144–148.
29. Kusiak, A.; Zheng, H.; Song, Z. On-line monitoring of power curves. *Renew. Energy* **2009**, *34*, 1487–1493. [[CrossRef](#)]
30. Lydia, M.; Kumar, S.S.; Selvakumar, A.I.; Kumar, G.E.P. Wind farm power prediction based on wind speed and power curve models. In *Intelligent and Efficient Electrical Systems*; Springer: Berlin/Heidelberg, Germany, 2018; pp. 15–24.
31. Ciulla, G.; D’Amico, A.; di Dio, V.; Brano, V.L. Modelling and analysis of real-world wind turbine power curves: Assessing deviations from nominal curve by neural networks. *Renew. Energy* **2019**, *140*, 477–492. [[CrossRef](#)]
32. Sheibat-Othman, N.; Othman, S.; Tayari, R.; Sakly, A.; Odgaard, P.F.; Larsen, L.F.S. Estimation of the wind turbine yaw error by support vector machines. *IFAC-PapersOnLine* **2015**, *48*, 339–344. [[CrossRef](#)]
33. Pandit, R.; Infield, D. Gaussian Process Operational Curves for Wind Turbine Condition Monitoring. *Energies* **2018**, *11*, 1631. [[CrossRef](#)]
34. Gill, S.; Stephen, B.; Galloway, S. Wind turbine condition assessment through power curve copula modelling. *IEEE Trans. Sustain. Energy* **2012**, *3*, 94–101. [[CrossRef](#)]
35. Vapnik, V.N. Principles of risk minimisation for learning theory. In *Advances in Neural Information Processing Systems*; Morgan Kaufman: San Mateo, CA, USA, 1992; pp. 831–838.
36. Vapnik, V.N. *Statistical Learning Theory*; John Wiley and Sons: New York, NY, USA, 1998.
37. Santos, P.; Villa, L.F.; Reñones, A.; Bustillo, A.; Maudes-Raedo, J.M. An SVM-based solution for fault detection in wind turbines. *Sensors* **2015**, *15*, 5627–5648. [[CrossRef](#)]
38. Dahhani, O.; El-Jouni, A.; Boumhidi, I. Assessment and control of wind turbine by support vector machines. *Sustain. Energy Technol. Assess.* **2018**, *27*, 167–179. [[CrossRef](#)]
39. Mohandes, M.A.; Halawani, T.O.; Rehman, S.; Hussain, A.A. Support vector machines for wind speed prediction. *Renew. Energy* **2004**, *29*, 939–947. [[CrossRef](#)]
40. Zeng, J.; Qiao, W. Short-Term Wind Power Prediction Using a Wavelet Support Vector Machine. *IEEE Trans. Sustain. Energy* **2012**, *3*, 255–264. [[CrossRef](#)]
41. Yan, H.; Mu, H.; Yi, X.; Yang, Y.; Chen, G. Fault Diagnosis of Wind Turbine Based on PCA and GSA-SVM. In Proceedings of the Prognostics and System Health Management Conference (PHM-Paris), Paris, France, 2–5 May 2019; pp. 13–17. [[CrossRef](#)]
42. Pandit, R.K.; Infield, D. SCADA based non-parametric models for condition monitoring of a wind turbine. *J. Eng.* **2019**, *2019*, 4723–4727.
43. Jin, T.; Tian, Z. Uncertainty analysis for wind energy production with dynamic power curves. In Proceedings of the 2010 IEEE 11th International Conference on Probabilistic Methods Applied to Power Systems, Singapore, 14–17 June 2010; pp. 745–750. [[CrossRef](#)]
44. De Brabanter, K.; De Brabanter, J.; Suykens, J.A.K.; De Moor, B. Approximate Confidence and Prediction Intervals for Least Squares Support Vector Regression. *IEEE Trans. Neural Netw.* **2010**, *22*, 110–120. [[CrossRef](#)]
45. Pandit, R.K.; Infield, D. SCADA-based wind turbine anomaly detection using Gaussian process models for wind turbine condition monitoring purposes. *IET Renew. Power Gener.* **2018**, *12*, 1249–1255. [[CrossRef](#)]
46. IEC Standard. *Wind Turbines—Part 12-1: Power Performance Measurements of Electricity Producing Wind Turbines (IEC 61400-12-1:2017)*; ICE: Geneva, Switzerland, 2017.
47. Pandit, R.K.; Infield, D.; Kolios, A. Gaussian process power curve models incorporating wind turbine operational variables. *Energy Rep.* **2020**, *6*, 1658–1669. [[CrossRef](#)]
48. Pandit, R.K.; Infield, D.; Carroll, J. Incorporating air density into a Gaussian process wind turbine power curve model for improving fitting accuracy. *Wind. Energy* **2019**, *22*, 302–315. [[CrossRef](#)]

49. Boser, B.E.; Guyon, I.M.; Vapnik, V. A training algorithm for optimal margin classifiers. In Proceedings of the Fifth Annual Workshop on Computational Learning Theory—COLT, Pittsburgh, PA, USA, 27–29 July 1992; p. 144, ISBN 089791497X. [[CrossRef](#)]
50. Rice, S.O. The distribution of the maxima of a random curve. *Am. J. Math.* **1939**, *61*, 409. [[CrossRef](#)]
51. Sun, J.; Loader, C.R. Simultaneous confidence bands for linear regression and smoothing. *Ann. Stat.* **1994**, *22*, 1328–1345. [[CrossRef](#)]
52. Eubank, R.L.; Speckman, P.L. Confidence Bands in Nonparametric Regression. *J. Am. Stat. Assoc.* **1993**, *88*, 1287. [[CrossRef](#)]
53. Abdi, H. Bonferroni and Sidak corrections for multiple comparisons. In *Encyclopedia of Measurement and Statistics*; Salkind, N.J., Ed.; Sage: Thousand Oaks, CA, USA, 2007; pp. 103–107.
54. Support Vector Machine Regression Models, Matlab Toolbox. Available online: <https://uk.mathworks.com/help/stats/support-vector-machine-regression.html> (accessed on 23 August 2020).
55. Tao, L.; Siqi, Q.; Zhang, Y.; Shi, H. Abnormal Detection of Wind Turbine Based on SCADA Data Mining. *J. Math. Probl. Eng.* **2019**, *2019*, 1–10. [[CrossRef](#)]

Publisher’s Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).