Bayesian Mixed Effect Model and Classical Time Series Forecasting on Wind Power

Kaiyan Xu

1 **Abstract:**

2     In this paper, I tried to model the natural factors that influence wind power and

3 conducted a forecast of wind power generation. Firstly, I built a Bayesian nonlinear model to

4 capture the general contribution of wind speed to wind power. Then I tried to add random

5 effects on hour and month to capture more variability and account for a lack of independence

6 in the original model. The second part of the paper is a classical time series model that

7 forecasts future wind power generation by using ARIMA method.

8

9 **1.  Introduction:**

10     The efficiency of using traditional energy has approached its threshold so keep relying on

11 traditional energy produces more cost than benefit. As a result, it is an imperative to gradually

12 turn to renewable energy like wind power. From 2015 to 2021, World renewable energy

13 investment increases from 14.8% in 2015, the year of Paris Agreement, to 19.8% in 2021

14 (IED). Data from IED also shows that net electricity produced by wind in the US increased

15 97.4% since 2015, in total 31378.5 GWh of electricity in 2021. Therefore, correctly modeling

16 and estimating wind power is becoming more important. Statistical models play a significant

17 role for windmill companies to schedule, dispatch and balance the demand and supply in the

18 electricity system. Accurate prediction minimizes the need for additional waste to reserve

19 energy and improve profit (Aoife et al. 2011).

20     However, capturing the variability of potential wind power generated is difficult. Wind is

21 weather dependent and highly stochastic. Wind is influenced by many other natural factors

22 including season, temperature, pressure, humidity, cloud cover, rainfall, etc. (Wang et al.

23 2011). Moreover, topographic characteristics like turbine position, turbine size, tower height,

24 elevation under different environments also influence the efficiency of generating power. In

25 general case, hybrid models with training data updated regularly have been found to be more

26   accurate model (Saroha et al. 2020). However, since in my dataset from a windmill in Texas

27   only has significant relationship with wind speed, it is sufficient to build a univariate model.

28       The second effect that leads to huge variability of wind power is the daily and seasonal

29   pattern. The hypothesis of this paper is that such daily and seasonally pattern have significant

30   underlying effects to the variability of wind power despite its high stochasticity. In my hourly

31   wind power data, I incorporate the hourly and monthly random effects to the original fixed

32   effect model and evaluated their effects. In the second part of the paper, the dataset is 10-year

33   monthly generated wind power data in the US. By including past seasonal pattern, the

34   accuracy of predicting future wind power is also assumed to have a huge improvement.

35

36   **2.   Materials and Methods**

37   *2.1 Wind Power Regression Analysis*

38       In the first part of the paper, hourly generated wind power along with other weather data

39   is measured in a windmill in Texas. By drawing correlation plot of wind power with each of

40   the four measurements: wind speed, wind direction, pressure, and temperature, only wind

41   speed has an obvious relationship with wind power, and the variability along with its curve is

42   small. Therefore, it is safe to only include wind speed as the regressor and determine there is

43   little omitted variable bias to do so. The shape of the curve between wind power and wind

44   speed is exponential. Therefore, I proposed a logistic shaped curve to simulate their

45   relationship. By calculating the correlation between wind speed and wind direction, pressure,

46   temperature, other weather measurements are uncorrelated with wind speed either. As a

47   result, there is almost no endogeneity in the model, i.e., E(residual|wind speed) = 0. There is

48   no need to add these variables as instruments to isolate the movements in wind speed that are

49   uncorrelated with model residuals, which in turn permits consistent estimation of coefficient

50   of wind speed (Stock and Francesco 2003).

51    The process model of the first fixed effect Bayesian model is:

52    $$\mu = a \times \frac{e^{b+c*x}}{1+e^{b+c*x}} ,$$

53    where a, b, c are parameters derived from the mean of posterior probability in the

54    Bayesian model. While logistic regression only constrains its range between [0,1], it is

55    required to have a multiplicator that magnify the range of logistic regression to roughly fit the

56    range of wind power, and parameter $a$ has such an effect. By implementing MCMC using 3

57    chains, the chains converge well so the results in posterior are valid. A plot of fitted curve

58    along with 95% confidence and prediction interval is shown in Figure 1. The fitted line fits

59    the observation well, and 95% confidence interval is quite close to the fitted line, which

60    means the fitted line is very accurate.

61    The data model of the Bayesian model is

62    $$y = dnorm(\mu, s),$$

63    where s is standard deviation, and it denotes the error in the model.   Here, it is assumed

64    that observations are normally distributed around the predicted value because in most cases,

65    the observed value follows normal distribution with the predicted curve. Although when wind

66    speed is smaller than 3m/s, no power is generated, and confidence interval and predictive

67    interval in low wind power might produces wind power that is lower than 0 in this case under

68    normal distribution, it is worthwhile to notice that this part of data only accounts for a very

69    rare situation and the distance lower the horizontal axis is very small. Therefore, a slight

70    negative predicted interval can be treated as 0 that there is no wind power at all.
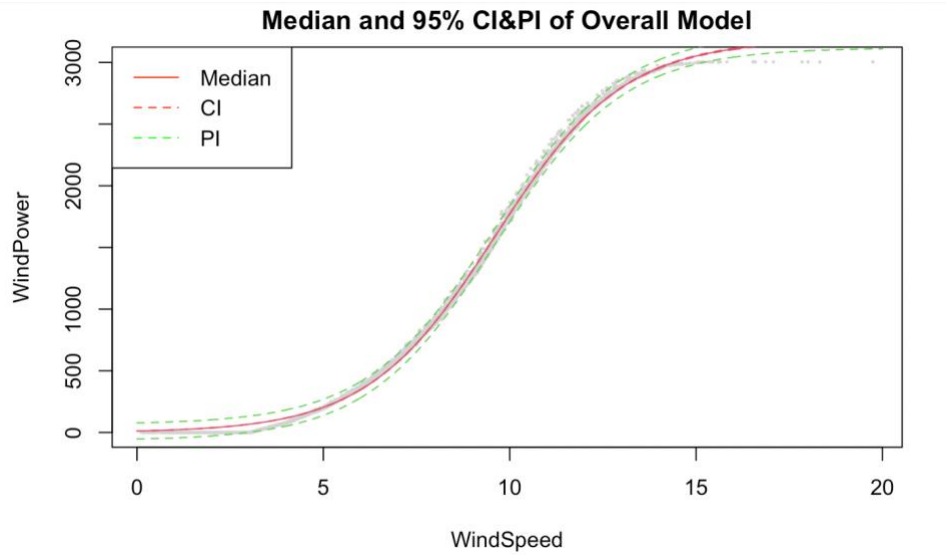
**Figure 1. Overall Bayesian Nonlinear Model**

Then, I divided the dataset to make sure if high wind power and low wind power have

significant different or opposite parameters. In such a case, I need to separate them into two

subgroups for further analysis. It is observed that wind power from 13:00 PM to 21:00 PM

and from January to June are much higher on average, so I set two subsets to run the previous

Bayesian model again. One subset consists of data from January to June at 13PM to 21PM;

another subset consists of data from July to December at 23PM to 11AM. After comparing

the parameters, I found that they are similar in two models, so there is no need to model high

wind power and low wind power separately.

The next step is to add hourly and seasonal random effect to the fixed effect model.

Random effects account for a lack of independence and unexplained variance within the same

group after repetition (Liu et al. 2016). Firstly, I added hour of the measurement as the first

random effect. Since it is a nonlinear model, every parameter can be set as random effect

(Barrowman et al. 2003). I added random effect on each of the three parameters a, b, c and on

the residual of the nonlinear model. The four process models become:

$$\mu_2 = a_i \times \frac{e^{b+c*x}}{1+e^{b+c*x}},$$

89
$$\mu_3 = a \times \frac{e^{b_i+c*x}}{1+e^{b_i+c*x}},$$

90
$$\mu_4 = a \times \frac{e^{b+c_i*x}}{1+e^{b+c_i*x}},$$

91
$$\mu_5 = a \times \frac{e^{b+c*x}}{1+e^{b+c*x}} + \alpha.h_i,$$

92    where $a_i$, $b_i$, $c_i$, $\alpha.h_i$ denotes that the hourly random effects added. In such an effect,

93    wind power measure in the same hour a day have the same parameter in the model. By

94    comparing DIC of the fixed effect model with the four random effect model above, it turns

95    out that setting residuals as random effect has lowest DIC (Table 1). Therefore, I considered

96    to keep adding random effects of month on the residual to see if it still improves DIC. The

97    process model becomes:

98
$$\mu_6 = a \times \frac{e^{b+c*x}}{1+e^{b+c*x}} + \alpha.h_i + \alpha.m_i,$$

99    where $\alpha.m_i$ denotes the monthly random effect added. Now the process model can

100   capture both the uncertainty from the same hour in a day and days within the same month.

101   This model with two random effects on residual turns out to have the lower DIC. In an

102   example, Figure 2 is a plot of best fitted line for January data shown in red. Still, confidence

103   interval is narrow so the uncertainty in estimate is small and it indicates that the best fitted

104   line is accurate. Most January data falls within the prediction interval, so the model captures

105   most variability of wind power in January.

106

| Process Model | DIC |
|---|---|
| $\mu = a \times \dfrac{e^{b+c*x}}{1 + e^{b+c*x}}$ | 43165 |
| $\mu_2 = a_i \times \dfrac{e^{b+c*x}}{1 + e^{b+c*x}}$ | 43133 |
| $\mu_3 = a \times \dfrac{e^{b_i+c*x}}{1 + e^{b_i+c*x}}$ | 43142 |

| | |
|---|---|
| $\mu_4 = a \times \dfrac{e^{b+c_i*x}}{1 + e^{b+c_i*x}}$ | 43134 |
| $\mu_5 = a \times \dfrac{e^{b+c*x}}{1 + e^{b+c*x}} + \alpha.h_i$ | 42373 |
| $\mu_6 = a \times \dfrac{e^{b+c*x}}{1 + e^{b+c*x}} + \alpha.h_i + \alpha.m_i$ | 41760 |

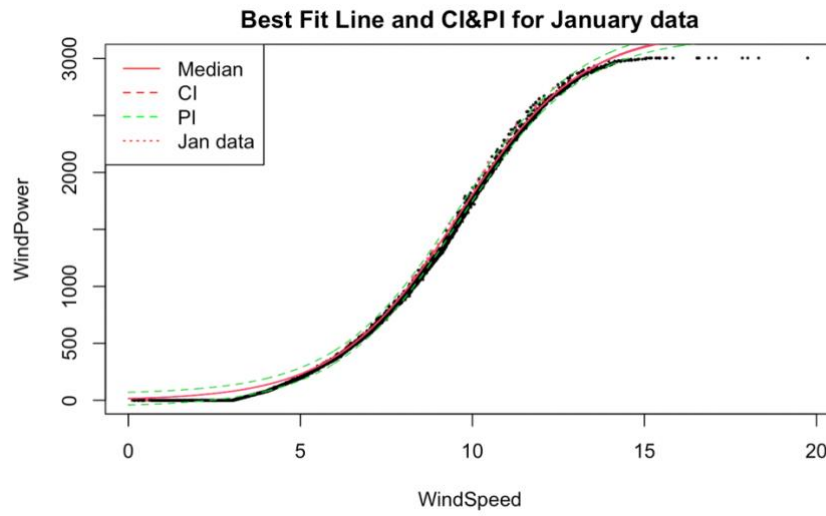107                               Table 1. DIC of Models

108



109

110                    Figure 2. Hourly and Monthly Random Effect Model January Data

111

112        There is still high autocorrelation over 100 lags between relevant measurements in the

113    residuals (Appendix a). Therefore, it is estimated that DIC will be much lower if accounting

114    for the remaining autocorrelation between closed measurements.

115

116    *2.2 Wind Power Forecast*

117        The second part of this paper is to forecast US wind power generation for the first 6

118    month of 2022 using classical ARIMA method in a frequentist context. The dataset for this

119    project consists of US monthly wind power generation from 2010-2021. It is extracted from

120    the larger dataset, "*Monthly Electricity Statistics - Monthly electricity production and trade*

121    *data for 47 countries*," from IEA.

122        The time series plot shows a clear upward trend and seasonal pattern in the data as

123    expected (Figure 3). Also, variance of time series is increasing so there are more variability of

124    wind power generation as average generation increases. Therefore, it is obvious that the

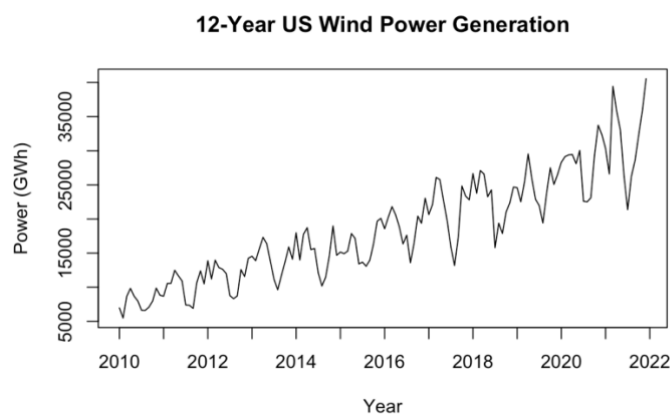125    process is nonstationary, and hence it requires transformation before catching the correct

126    model.



127

128                    Figure 3. Time Series Plot of US Monthly Wind Power Generation

129

130        To transform the process to be stationary, I firstly applied Box-Cox transformation to

131    stabilize the variance. I applied order 1 differencing to eliminate the trend and order 12

132    differencing to eliminate the seasonality as it is monthly data. Starting from here, I can fit

133    models based on the ACF and PACF plot of residuals (Appendix b).

134        The method is to fit a model based on autocorrelation and partial autocorrelation. From

135    the ACF and PACF plot, for nonseasonal part, ACF cuts off after lag 1 and PACF can be

136    treated as either decaying to 0 or cut off after lag 4. For seasonal part, ACF either cuts off

137    after seasonal lag 1 (lag12) or decays to zero, the same as PACF. To find the model with

138     lowest $AIC_C$, I fit 6 relevant models derived from the correlation plots. The table (Appendix

139     c) shows the 6 models and their $AIC_C$.

140      Based on the Table 1, $SARIMA(4,1,1) * (1,1,0)_{12}$ and $SARIMA(0,1,1) * (1,1,0)_{12}$

141     have lowest $AIC_C$. However, a diagnostic show that the normality assumption of fitting an

142     ARIMA model is violated in these two models, so they are inadequate models. The model

143     with third lowest $AIC_C$, on the other hand, shows normality on its residuals and that its

144     standardized residuals follow white noise process (Appendix d). Therefore, it is an adequate

145     and relative accurate model for prediction. The result of prediction is shown in Figure 5. By

146     taking out the last 6-month data as testing set and refit the data using the first 138 data as

147     training set, all 6 testing data is contained in the 80% confidence interval. Therefore, the

148     ARIMA model here is accurate in predicting the future wind power generation.
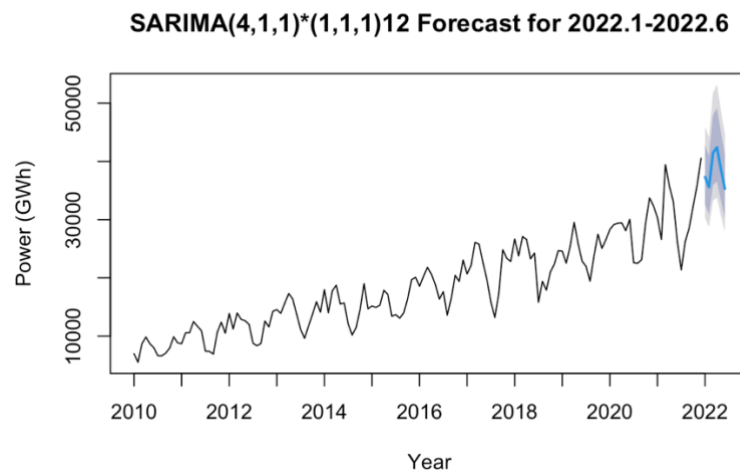


SARIMA(4,1,1)*(1,1,1)12 Forecast for 2022.1-2022.6

149

150                          Figure 5. $SARIMA(4,1,1) * (1,1,1)_{12}$ Forecast

151

152     **3.   Results**

153      In Texas Windmill Dataset, the model with the lowest DIC is the mixed model with hour

154     and month as random effects. It verifies that there is hourly and seasonal pattern in wind

155     power. Within the same hour a day and same month, the variability can be captured by the

156   same random effect. When modeling seemingly highly stochastic data, considering random

157   effects improves the accuracy of the model by borrowing strength across the dataset.

158       In US wind power generation dataset, the forecasting wind power from this model does

159   not deviate from the historical pattern a lot. It is on a high position and follows past average

160   seasonal pattern: from December to February, wind power decreases; from March to April,

161   wind power is high; from May to June, wind power goes down again. Therefore, seasonal

162   pattern is significant. This forecast assumes that no intervention would occur. Based on the

163   forecast, windmill companies can plan their production and storage of wind power

164   accordingly, and government can manage production of other renewable energy when during

165   the months when wind power generation is low.

166

167   **4.  Discussion**

168       There are many constraints in this paper. Firstly, the Texas Windmill Dataset only has

169   wind speed that has obvious relationship with wind power, which is not representative of

170   most wind power generation pattern in real life. It is necessary to find another dataset with

171   more useful covariates and build a more complex model to comprehensively understand

172   effects of natural factors on wind power.

173       Secondly, I did not try enough combination of random effects on nonlinear model

174   parameters. I only assumed one of the parameters are random effects. When adding the

175   second random effect, month, I only added to the best model after adding the first random

176   effect. Therefore, it is possible that I missed the model with lowest DIC. To solve this I need

177   to use greedy method to compare every combination of parameters.

178       Thirdly, I did not consider the remaining autocorrelation after adding the random effects.

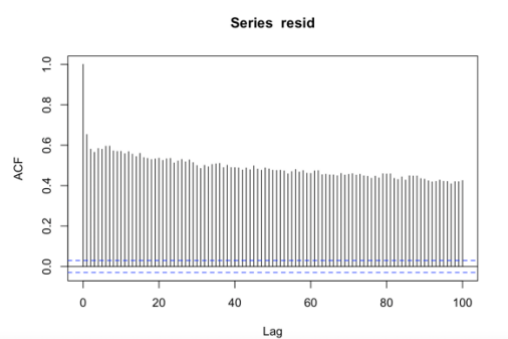179   There is still high autocorrelation between relevant data points over 100 lags. I need to add a

180 covariance matrix as the variance of data model to account for different covariance at each

181 point.

182

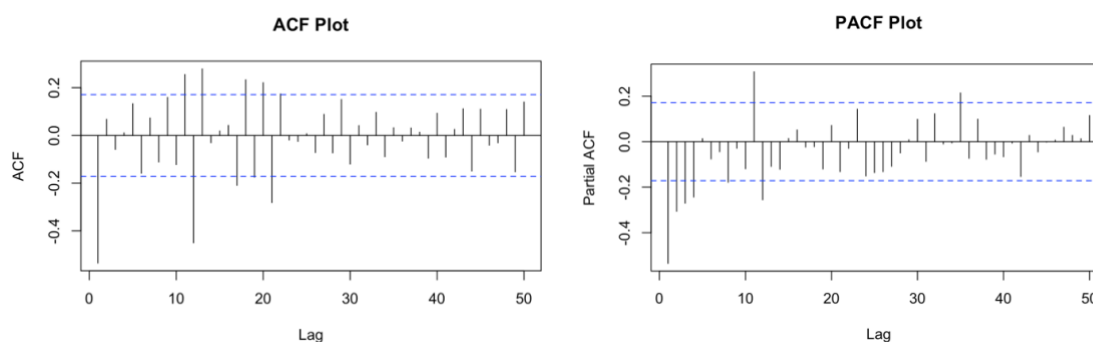183 **5. Acknowledgement**

184 I thank Professor Michael Dietze for his instruction and help for the whole semester.

185

186 **6. Appendix**

187 a. Autocorrelation between Measurements



188

189 b. ACF and PACF of Differenced Time Series



190

191 c. $AIC_C$ of ARIMA Models

| Model | $AIC_C$ |
|---|---|
| SARIMA(0,1,1)(1,1,1)[12] | -198.31 |
| SARIMA(0,1,1)(0,1,1)[12] | -200.43 |

| SARIMA(0,1,1)(1,1,0)[12] | -174.25 | 192 |
| SARIMA(4,1,1)(0,1,1)[12] | -192.09 | 193 |
| SARIMA(4,1,1)(1,1,0)[12] | -166.68 | 194 |
| SARIMA(4,1,1)(1,1,1)[12] | -189.88 | 195 |

196

197    d.  Diagnostics of ARIMA Model



Normal Q-Q Plot for ARIMA(4,1,1)(1,1,1)[12]

198



199

200

201    **7.  References**

202    World Energy Investment 2021 Datafile, IED, https://www.iea.org/data-and-statistics/data-

203         product/world-energy-investment-2021-datafile

204    Monthly Electricity Statistics, IED,

205  https://www.iea.org/data-and-statistics/data-product/monthly-electricity-statistics

206  Aoife M. Foleyabd, Paul G. Leahyab, Antonino Marvugliac, Eamon J.McKeoghab. Current

207    methods and advances in forecasting of wind power generation. 2011.

208    https://doi.org/10.1016/j.renene.2011.05.033

209  Xiaochen Wang, Peng Guo, Xiaobin Huang, A Review of Wind Power Forecasting Models,

210    Energy Procedia, Volume 12, 2011, Pages 770-778, ISSN 1876-6102,

211    https://doi.org/10.1016/j.egypro.2011.10.103.

212  Sumit Saroha, Sanjeev Kumar Aggarwal and Preeti Rana. Wind Power Forecasting, 2020,

213    DOI: 10.5772/intechopen.94550. https://www.intechopen.com/chapters/74076

214  Stock, James H., and Francesco Trebbi. 2003. "Who Invented Instrumental Variable

215    Regression?" Journal of Economic Perspectives 17: 177–194.

216  Liu, Xiaolei et al. "Iterative Usage of Fixed and Random Effect Models for Powerful and

217    Efficient Genome-Wide Association Studies." PLoS genetics vol. 12,2 e1005767. 1 Feb.

218    2016, doi:10.1371/journal.pgen.1005767

219  Barrowman, Nicholas J., et al. "The Variability among Populations of Coho Salmon in the

220    Maximum Reproductive Rate and Depensation." Ecological Applications, vol. 13, no. 3,

221    2003, pp. 784–93, http://www.jstor.org/stable/4134695. Accessed 7 May 2022.