1 **A robust error correction method for numerical weather**
2 **prediction wind speed based on Bayesian optimization,**
3 **Variational Mode Decomposition, Principal Component**
4 **Analysis, and Random Forest: VMD-PCA-RF (version**
5 **1.0.0)**

6 Shaohui Zhou[1], Chloe Yuchao Gao[2*], Zexia Duan[1], Xingya Xi[3], and Yubin Li[1]
7 [1]Collaborative Innovation Centre on Forecast and Evaluation of Meteorological Disasters, Key
8 Laboratory for Aerosol-Cloud-Precipitation of China Meteorological Administration, School of
9 Atmospheric Physics, Nanjing University of Information Science and Technology, Nanjing, 210044,
10 China.
11 [2]Department of Atmospheric and Oceanic Sciences and Institute of Atmospheric Sciences, Fudan
12 University, Shanghai, 200438, China.
13 [3]School of Atmospheric Sciences, Sun Yat-sen University, and Southern Marine Science and
14 Engineering Guangdong Laboratory (Zhuhai), Zhuhai, 519082, China

15 *Correspondence to*: gyc@fudan.edu.cn

16 **Abstract.** Accurate wind speed prediction is crucial for the safe utilization of wind resources. However,

17 current single-value deterministic numerical weather prediction methods employed by wind farms do

18 not adequately meet the actual needs of power grid dispatching. In this study, we propose a new hybrid

19 forecasting method for correcting 10-meter wind speed predictions made by the Weather Research and

20 Forecasting (WRF) model. Our approach incorporates Variational Mode Decomposition (VMD),

21 Principal Component Analysis (PCA), and five artificial intelligence algorithms: Deep Belief Network

22 (DBN), Multilayer Perceptron (MLP), Random Forest (RF), eXtreme Gradient Boosting (XGBoost),

23 light Gradient Boosting Machine (lightGBM), and the Bayesian Optimization Algorithm (BOA). We

24 first construct WRF-predicted wind speeds using the Global Prediction System (GFS) model output

25 based on prediction results. We then perform two sets of experiments with different input factors and

26 apply BOA optimization to debug the four artificial intelligence models, ultimately building the final

27 models. Furthermore, we compare the forementioned five optimal artificial intelligence models suitable

28 for five provinces in southern China in the wintertime: VMD-PCA-RF in December 2021 and

29 VMD-PCA-lightGBM in January 2022. We find that the VMD-PCA-RF evaluation indexes exhibit

30 relative stability over nearly a year: correlation coefficient (R) is above 0.6, accuracy rate (FA) is above

31 85 %, mean absolute error (MAE) is below 0.6 m/s, root mean square error (RMSE) is below 0.8 m/s,

32    relative mean absolute error (rMAE) is below 60 %, and relative root mean square error (rRMSE) is

33    below 75 %. Thus, for its promising performance and excellent year-round robustness, we recommend

34    adopting the proposed VMD-PCA-RF method for improved wind speed prediction in models.


35    **1 Introduction**

36        Sustainable energy plays a vital role in reducing carbon footprint and increasing system reliability

37    (Hanifi et al., 2020). As renewable energy sources have a negligible carbon footprint, they have

38    become the preferred choice for many industries in the power sector (Dhiman and Deb, 2020). Among

39    these sources, wind energy is a crucial low-carbon energy technology with the potential to become a

40    sustainable energy source (Tascikaraoglu and Uzunoglu, 2014). In 2022, the global wind power

41    capacity reached 906 GW, with a 9 % year-on-year increase due to a newly installed capacity of 77.6

42    GW. The global onshore wind market increased by 68.8 GW, while facing a 5 % growth decline

43    compared to the previous year. Such change is attributed to a slowdown in China and the U.S., the

44    world's two largest wind markets that account for over two-thirds of the world's onshore wind farm

45    installations (Joyce and Feng, 2023). Therefore, accurate and stable wind speed prediction (WSP) is

46    very important for the safe and stable operation of the power grid system and improving the utilization

47    rate of wind energy and economic development (Guo et al., 2021; Xiong et al., 2022; Tang et al.,

48    2021).

49        Current WSP algorithms are primarily categorized into physical algorithms (Zhao et al., 2016),

50    statistical algorithms (Wang and Hu, 2015; Barthelmie et al., 1993), machine learning (ML) algorithms

51    (Huang et al., 2019; Salcedo-Sanz et al., 2011; Ma et al., 2020), and hybrid algorithms (Deng et al.,

52    2020; Xu et al., 2021; Zhao et al., 2019; Xiong et al., 2022; Tang et al., 2021). Physical methods, such

53    as numerical weather prediction (NWP), are commonly used in wind speed forecasting. NWP, which

54    accounts for atmospheric processes and physical laws, solves discrete mass, momentum, and energy

55    conservation equations along with other fundamental physical principles, establishing itself as a widely

56    adopted and reliable physical method. Currently, the High-resolution Limited Area Model (HIRLAM)

57    (Landberg, 1999), the European Center for Medium-Range Weather Forecast (ECMWF) model, the

58    fifth-generation mesoscale model (MM5) (Salcedo-Sanz et al., 2009), and the Weather Research and

59    Forecasting Model (WRF) (Prósper et al., 2019) are extensively utilized for wind speed prediction.

60    However, NWP modeling faces challenges due to the selection of parameterization schemes, such as

61    model microphysics and systematic errors, which exhibit temporal and spatial differences and

62    uncertainties. These uncertainties hinder the accuracy of NWP models in wind speed prediction,

63    making it difficult to meet the rising demands of the grid system (Zhao et al., 2019; Xu et al., 2021).

64    Studies have demonstrated that enhancing the accuracy of numerical weather prediction (NWP)

65    models and correcting prediction errors can effectively minimize the errors associated with wind speed

66    prediction. These research endeavors have typically sought to optimize the physical and dynamic

67    parameters of the NWP model (Cheng et al., 2013), refine the model structure (Jiménez and Dudhia,

68    2012), or improve the accuracy of model inputs through preprocessing and denoising techniques (Xu et

69    al., 2015). Additionally, improving initial field error through methods, such as target observation and

70    data assimilation (Williams et al., 2013), can also minimize wind speed errors predicted by NWP

71    models.

72    Physical methods are generally more appropriate for long-term wind speed prediction, such as

73    those 48-72 hours in advance, while their practical application in short-term forecasting is limited

74    (Zhao et al., 2019; Deng et al., 2020; James et al., 2018). In contrast, statistical methods utilize

75    historical data to establish a relationship between input and output variables and are therefore

76    well-suited for short-term wind speed prediction. They are usually time series models, such as

77    Autoregressive Moving Average (ARMA) (Erdem and Shi, 2011) and Autoregressive Integrated

78    Moving Average (ARIMA) (Wang and Hu, 2015). Whereas filtering models (Cassola and Burlando,

79    2012; Chen and Yu, 2014), machine learning models (Hu et al., 2013), and hybrid models (Huang et al.,

80    2019) have been gradually developed to further improve wind speed prediction accuracy.

81    With purely statistical models becoming less suitable for wind speed predictions beyond 6 hours,

82    the use of a combination of physical and statistical methods has gained growing interest (Zjavka, 2015;

83    Xu et al., 2021). The error correction model improves the accuracy of the NWP model by training the

84    relationship between the NWP predictor variables and the observed correlation variables (Sun et al.,

85    2019). However, traditional error prediction models rely solely on historical wind speed sequences as

86    input factors (Deng et al., 2020; Guo et al., 2021) and do not incorporate the characteristic

87    meteorological factors forecasted by the WRF model. Studies have shown that considering all relevant

88    historical meteorological factors can lead to more accurate predictions compared to only taking into

89    account historical wind speed (Zhang et al., 2019c). Therefore, it is crucial to include meteorological

90    characteristic factors as input in the prediction model.

91        For an error prediction model, wind speed is the most important input factor. Traditionally, the

92    error prediction model uses historical wind speed data as input, without any feature selection. Feature

93    selection methods, such as filtering methods, are commonly used in time series analysis. Currently,

94    empirical mode decomposition (EMD) (Liu et al., 2018; Guo et al., 2012), ensemble empirical mode

95    decomposition (EEMD) (Wang et al., 2017), wavelet decomposition (WD) (Zhang et al., 2019b),

96    variational mode decomposition (VMD) (Hu et al., 2021; Zhang et al., 2019a), and other filtering

97    methods are used to select key features in the wind speed data. As mentioned above, studies have

98    shown that these feature selection methods can effectively extract the hidden features in the wind speed

99    series to improve wind speed prediction accuracy. However, despite the effectiveness of wind speed

100    filtering methods in wind speed prediction, only a few studies have applied these methods to the

101    correction of wind speed errors in NWP forecasting (Xu et al., 2021; Li et al., 2022).

102        In addition, traditional error correction methods generally adopt linear regression (Dong et al.,

103    2013), multiple linear regression (Liu et al., 2016), machine learning (Salcedo-Sanz et al., 2011), and

104    deep learning algorithms (Zhang et al., 2019c). However, the efficacy of machine learning and deep

105    learning algorithms is highly dependent on the selection of model parameters (Guo et al., 2021; Xiong

106    et al., 2022). The Bayesian optimization algorithm (Li and Shi, 2010; Guo et al., 2021) is considered a

107    relatively advanced algorithm for optimizing model parameters and has been widely used in MATLAB

108    and Python packages.

109        In this study, we investigate a multi-step wind speed forecasting model that combines NWP

110    simulation and an error correction strategy. We present two sets of experiments divided into three steps:

111    (1) we use the first group of experiments to extract hidden features from various meteorological

112    elements forecasted by NWP; The second group of experiments mainly focuses on the wind speed

113    forecast of NWP, and the VMD-PCA algorithm is used to extract the hidden features in the forecasted

114    wind speed; each set of experimental input factors is matched with the actual 10-meter wind speed data

115    of 410 stations in time and space; (2) we employ four advanced machine learning algorithms optimized

116    by the BOA algorithm, and DBN deep learning algorithm to train the two groups of experiments and

117    perform 5-fold cross-validation; and (3) we analyze six distinct wind speed error indicators to compare

118     and identify the most suitable wind speed error correction schemes for the five southern provinces in

119     winter and throughout most of the year. The remainder of this paper is organized into sections

120     discussing the effects of the BOA-VMD-PCA approach, the interpretability of RF feature importance,

121     and the stability analysis of the proposed models.

**2 Data and methods**

**2.1 Data**

124         The target observation data includes 2-m air temperature, 2-m specific humidity, 10-meter wind

125     speed, surface pressure, and precipitation. These data are collected on equivalent latitude and longitude

126     grid scale, primarily from five provinces in China: Guangdong, Guangxi, Yunnan, Guizhou, and

127     Hainan, covering a geographical range of 15-32.97°N and 94-120.97°E. The spatial resolution of the

128     grid is 0.03° × 0.03° and the temporal resolution is 1 hour. The dataset is constructed through the

129     integration of multiple sources, including ground and satellite data, and is refined using advanced

130     techniques such as multi-grid variational assimilation, physical inversion, and terrain correction. This

131     dataset exhibits superior quality in comparison to other products, offering higher spatial and temporal

132     resolutions. For the purposes of this paper, the 10-meter wind speed data is interpolated across 410
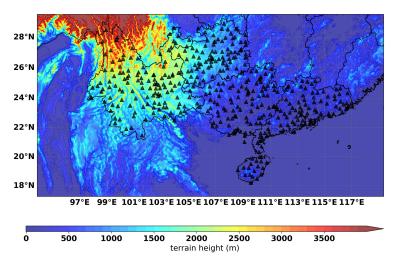
133     sites, as illustrated in Figure 1.



**Figure 1: The elevation map of the five southern provinces in china (black triangles represent weather stations).**

137

## 2.2 Methods

### 2.2.1 WRF simulation

140    The WRF 4.2 model, developed by the United States' National Center for Environmental
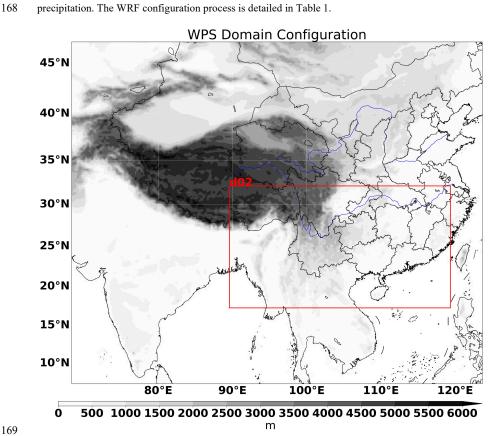
141 Prediction (NCEP), represents a new generation of mesoscale numerical models with numerous

142 applications in research forecasting. The WRF model, written in Fortran 90 language, offers

143 advantages such as portability, scalability, and high efficiency. It employs Arakawa C-grid points in the

144 horizontal direction and terrain-following mass coordinates in the vertical direction. When forecasting

145 meteorological elements, the WRF model uses the US Global Weather Forecast Data (GFS) developed

146 by NCEP and the National Center for Atmospheric Research (NCAR). The GFS system includes data

147 related to the atmosphere and land variables, such as temperature, precipitation, and wind data. The

148 system is updated every 6 hours, at 0:00, 6:00, 12:00, and 18:00 UTC, and provides predictions for the

149 subsequent eight days. Given that the time scale of the meteorological station data in the study area is 1

150 hour, the forecast data time interval of the WRF model is also set to 1 hour. As a widely used

151 numerical weather forecast model, the WRF model is suitable for weather studies from a few meters to

152 several thousand kilometers. Therefore, this paper uses the WRF model to predict 10-meter wind speed

153 as the input factor for the error correction model (Xu et al., 2021).

154    Using the WRF model in combination with daily data resolution of $0.25° \times 0.25°$, the model

155 initiates at 18:00 UTC and generates forecasts every 3 hours for a total duration of 102 hours. The

156 regular Global Forecast System (GFS) forecast field data serve as the initial field and lateral boundary

157 conditions for the WRF model. Surface static data, such as terrain, soil data, and vegetation coverage,

158 are derived from the Moderate Resolution Imaging Spectroradiometer (MODIS) satellite with a

159 resolution of 15 seconds (approximately 500 meters). Incorporating a two-layer grid nesting

160 configuration, the forecast area is illustrated in Figure 2. The grid dimensions are 600×500 and

161 967×535, with horizontal grid resolutions of 9 km and 3 km, respectively. The grid center points are set

162 at 29°N and 96°E. The "CONUS" parameterization scheme is used, including the Thompson

163 microphysics scheme, the Tiedtke cumulus parameterization scheme, the RRTMG long-wave and

164 short-wave radiation schemes, the Mellor-Yamada-Janjić (MYJ) boundary layer and near-surface

165 parameterization schemes, and the MYJ surface layer parameterization scheme. The Noah Land

166     Surface Model (LSM) is utilized for the surface process plan, generating a WRFOUT numerical

167     weather forecast file including meteorological elements such as temperature, humidity, and

168     precipitation. The WRF configuration process is detailed in Table 1.



169

170     **Figure 2: Schematic diagram of the simulation area of the WRF model.**

171

172                       **Table 1: WRF configuration scheme**

| Model (Version) | WRF (V4.2) | |
|---|---|---|
| Domains | D1 | D2 |
| Horizontal grid points | 600*500 | 967*535 |
| Δx (km) | 9 | 3 |
| Vertical layers | 58 | |
| Longwave radiation | RRTMG (Iacono et al. 2008) | |
| Shortwave radiation | RRTMG (Iacono et al. 2008) | |
| Land surface | Noah LSM (Chen et al. 1997) | |
| Surface layer | MYJ (Janjic 1994) | |
| Microphysics | Thompson (Thompson et al. 2008) | |

| Boundary layer | MYJ (Janjic 1994) |
|---|---|
| Cumulus | Tiedtke (Tiedtke 1989, Zhang et al. 2011) |

173

### 2.2.2 Variational mode decomposition

As a new filtering method, VMD is robust in feature selection. The VMD algorithm decomposes a time series signal into several intrinsic mode functions (Isham et al., 2018). The sum of the modes equals the original signal, and the sum of the bandwidths is the smallest. The analysis signal is calculated using the Hilbert transform to estimate the modal bandwidth. The optimization model is described as

$$\left\{ \min_{\{u_k\},\{\omega_k\}} \left\{ \sum_{k=1}^{K} \| \partial_t \left[ \left( \delta(t) + \frac{j}{\pi t} \right) u_k(t) \right] e^{-j\omega_k t} \Big\|_2^2 \right\} \; s.t. \quad \sum_{k=1}^{K} u_k = v \right. \qquad (1.1)$$

where $K$ is the total number of modes, $u_k$ is the decomposed $K$-th mode, $w_k$ is the corresponding center frequency, and $v$ is the time-series signal, representing the wind speed sequence predicted by the WRF model in this study.

The above constrained problem can be transformed into an unconstrained problem using the Lagrangian function:

$$L\left(\{u_k\},\{\omega_k\},\lambda\right) = \omega_k^{n+1} = \frac{\int_0^{\infty} \omega \left| \hat{u}_k(\omega) \right|^2 d\omega}{\int_0^{\infty} \left| \hat{u}_k(\omega) \right|^2 d\omega} \sum_{k=1}^{K} \| \partial_t \left[ \left( \delta(t) + \frac{j}{\pi t} \right) u_k(t) \right.$$

$$\times ] e^{-j\omega_k t} \|_2^2 + \| v(t) - \sum_{k=1}^{K} u_k(t) \|_2^2 + \left\langle \lambda(t), v(t) - \sum_{k=1}^{K} u_k(t) \right\rangle \qquad (1.2)$$

where $\alpha$ is the penalty parameter and $\lambda(t)$ is the Lagrange multiplier.

Then we update $u_k$, $w_k$, and $\lambda$ using the alternating direction method of the multiplier:

$$\hat{u}_k^{n+1}(\omega) = \frac{\hat{v}(\omega) - \sum_{i \neq k} \hat{u}_i(\omega) + \frac{\hat{\lambda}(\omega)}{2}}{1 + 2\alpha \left( \omega - \omega_k \right)^2} \qquad (1.3)$$

$$\omega_k^{n+1} = \frac{\int_0^{\infty} \omega \left| \hat{u}_k(\omega) \right|^2 d\omega}{\int_0^{\infty} \left| \hat{u}_k(\omega) \right|^2 d\omega} \qquad (1.4)$$

$$\hat{\lambda}^{n+1}(\omega) = \hat{\lambda}^n(\omega) + \tau \left[ \hat{v}(\omega) - \sum_{k=1}^{K} \hat{u}_k^{n+1}(\omega) \right] \qquad (1.5)$$

194    where $\tau$ is the update parameter.

195         When the accuracy (left side of the following expression) meets the following condition, $u_k$, $w_k$

196    and $\lambda$ would stop updating:

197
$$\sum_{k=1}^{K} \frac{\parallel \hat{u}_k^{n+1} - \hat{u}_k^n \parallel_2^2}{\parallel \hat{u}_k^n \parallel_2^2} < \varepsilon \qquad (1.6)$$

198    where $\varepsilon$ is the tolerance of the convergence criterion.

199         The VMD algorithm is implemented to decompose the wind speed signal predicted by the WRF

200    model. When using multiple sub-signals instead of the original signal, more features of the wind speed

201    can be obtained. Therefore, it is beneficial to improve the prediction accuracy when using the

202    sub-signal as input to the error correction model (Xu et al., 2021; Li et al., 2022).

203    **2.2.3 Principal Component Analysis**

204         Subsequences obtained by VMD usually have several illusory components. Using PCA to extract

205    the principal components of subsequences increases the number of features input to the model and

206    reduces the dimension of the data decomposed by VMD. When pcs are used as the input of the error

207    prediction algorithm, the pcs fully reflect the characteristics of the subsequence and reduce the model

208    complexity. The pcs $y_k$, $k=1, 2, …, K$ of the subsequence matrix U and the cumulative contribution rate

209    $\eta_n$ of first $n$ principal components are expressed as:

210
$$y_k = c_k' U \qquad (1.7)$$

211
$$\eta_n = \frac{\sum_{k=1}^{n} \lambda_k}{\sum_{k=1}^{K} \lambda_k} \qquad (1.8)$$

212    where $c_k$ is the corresponding characteristic unit vector, with $k=1, 2, …, K$;  $\lambda_k$  is the characteristic

213    root, with  $\lambda_1 \geq \lambda_2 \geq … \geq \lambda_K$ .

214    **2.2.4 Proposed hybrid forecasting algorithms**

215         This study used five machine learning algorithms to conduct ten experiments across two main

216    paths. The first path involves increasing the variables related to wind speed in the forecast field, while

217    the second path entails extracting potential characteristic information of the forecast wind speed

218    through VMD and PCA and reducing the characteristic quantity of other forecast data. The overarching

219    goal is to achieve accurate correction of the forecast field wind speed. The flowchart of the artificial

220    intelligence models used to correct the WRF predicted wind speed for the two main experimental paths

221    is illustrated in Figure 3 and comprises the following three steps:

222    Step 1. Data fusion, cleaning, and standardization: As depicted in Figure 3, this paper proposes

223    two distinct experimental paths, with the primary difference being the selection of input variables. In

224    Experiment 1, as shown in Figure 6(c), 12 sets of data are selected from the WRF forecast field,

225    including altitude, 10-meter wind speed, latitude, longitude, surface pressure, relative humidity,

226    10-meter meridional wind, 10-meter zonal wind, 2-meter temperature, 2-meter dew point temperature,

227    10-meter wind direction, and hourly precipitation. Experiment 2, as illustrated in Figure 6(d), derives 8

228    sets of data by reducing the selected WRF field forecast data, including altitude, 10-meter wind speed,

229    latitude, longitude, surface pressure, relative humidity, 2-meter temperature, and hourly precipitation.

230    The focus is on unearthing hidden characteristic information of forecast wind speed. In this experiment,

231    the wind speed is decomposed into 9 Intrinsic Mode Functions (IMF) using VMD. Subsequently, a

232    low-dimensional wind speed vector is extracted from the 9 IMF components via PCA dimensionality

233    reduction, and all data are concatenated to construct the input factors for the model in Experiment 2.

234    Missing and outlier values are removed from the dataset. The two experiments standardize 12 sets of

235    meteorological elements (8 sets of meteorological elements in Figure 4, 9 IMF components, and three

236    PCA vectors in Figure 5) and wind speed observation data, respectively. Standardization addresses the

237    issue of varying meteorological factor values during training, which may result in different

238    contributions. In this paper, the 24-hour forecast data correspond to the observation data of the

239    subsequent 24 hours. The dataset spans from 00:00 on December 1, 2021, to 23:00 on February 28,

240    2022, totaling 2160 hours and encompassing 410 weather stations. Consequently, the original dataset

241    comprises 2160*410 samples, with each sample containing 12 meteorological features in Experiment 1

242    and 20 input features in Experiment 2.

243    Step 2. BOA optimization of AI models and cross-validation: In this study, the dataset is

244    partitioned into training, validation, and test sets in accordance with the time series. February 2022

245    serves as the training and validation sets, while December 2021 and January 2022 constitute the test set.

246    The training and validation sets are divided based on five-fold cross-validation. Both experiments

247    employ five machine learning algorithms (DBN, MLP, RF, XGBoost, and LightGBM) to construct

248    distinct machine learning models. Concurrently, this paper utilizes the BOA algorithm to tune the

249    parameters of all models, except for DBN, resulting in the optimal hyperparameters for each model.

250        Step 3. Model evaluation and error analysis: The trained machine learning models are applied to

251    the test set to obtain the revised wind speed data, and ultimately, the accuracy of all models is assessed

252    through the wind speed evaluation index. The ultimate goal here is to identify the best wind speed

253    correction model suitable for the entire year. Accordingly, the generalization of all models is evaluated

254    across other seasonal months of the year, culminating in the selection of the best model.
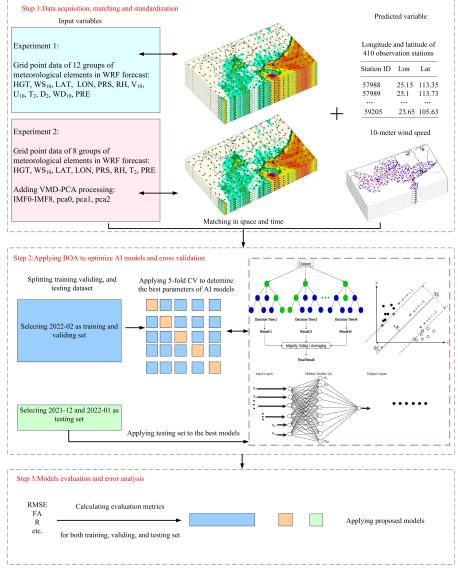


255

256

**Figure 3: Flowchart of the AI model used to correct WRF-predicted wind speeds in the two main experimental pathways.**
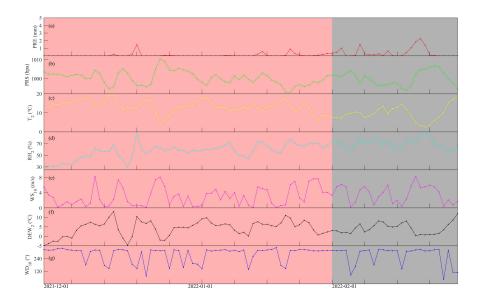
259



260

**Figure 4: Daily average hourly rainfall (a), surface pressure (b), 2-meter temperature (c), 2-meter relative humidity (d), 10-meter wind speed (e), 2-meter dew point temperature (f), and 10-meter wind direction ( g) which are located at Guangdong Lechang Station from December 1, 2021, to February 28, 2022. (February 2022 represents the training and verification sets, and December 2021 to January 2022 represents the testing set).**

266

267



268    **Figure 5: Three-dimensional view of 12 wind speed components after VMD and PCA processing of the**
269    **10-meter forecast wind speed at Lechang Station in Guangdong from December 1, 2021, to February 28,**
270    **2022.**

271

272    **2.2.5 Evaluation indicators**

273    There are many commonly used predictive effect evaluation indicators. This article uses the

274    following evaluation indicators: correlation coefficient (R), root mean square error (RMSE), mean

275    absolute error (MAE), relative root mean square error (rRMSE), relative mean absolute error (rMAE),

276    percentage of absolute error not greater than 1 m/s (FA). Six error indicators are used to evaluate the

277    correction results of short-term wind speed forecasts of wind farms. The formula for calculating the

278    error index is as follows:

279
$$R = \frac{\sum_{i}^{n}(y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}\sqrt{\sum_{i=1}^{n}(\hat{y}_i - \bar{\hat{y}})^2}} \qquad (1.9)$$

280
$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)^2} \qquad (1.10)$$

281
$$MAE = \frac{1}{n}\sum_{i=1}^{n}|\hat{y}_i - y_i| \qquad (1.11)$$

13

282
$$rRMSE = \left[ \sqrt{\frac{1}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)^2} \Big/ \left(\frac{1}{n}\sum_{i=1}^{n} y_i\right) \right] \times 100\% \qquad (1.12)$$

283
$$rMAE = \left(\frac{1}{n}\sum_{i=1}^{n}|\hat{y}_i - y_i| \Big/ \left(\frac{1}{n}\sum_{i=1}^{n} y_i\right)\right) \times 100\% \qquad (1.13)$$

284
$$FA = N_r / N_f \qquad (1.14)$$

285 Among them, $n$ represents the number of samples, $\hat{y}_i$ represents the $i$-th predicted value, $y_i$

286 represents the $i$-th actual value; $N_r$ represents the number of wind speed absolute errors not greater than

287 1 m/s, and $N_f$ represents the number of research samples.

288

289 **3 Results**

290 **3.1 Experiment 1 evaluation**

291 In Experiment 1, the BOA optimization algorithm was applied to five AI models to correct the

292 10-meter wind speed forecasted by WRF. There were 12 meteorological element features to establish

293 five different AI models (see Table 2 for the hyper-parameters of the five AI models). The training,

294 validation, and testing results for 10-meter wind speed are shown in Figures S1-5 in the supplementary

295 material. The RMSE values between the predicted and the observed value of the training set (validation

296 set) in the lightGBM, XGBoost, RF, DBN, and MLP models are 0.41 m/s (0.54 m/s), 0.31 m/s (0.56

297 m/s), 0.52 m/s (0.57 m/s), 0.59 m/s (0.62 m/s) and 0.73 m/s (0.73 m/s). The FA are 0.98 (0.94), 0.99

298 (0.93), 0.94 (0.93), 0.92 (0.91), and 0.88 (0.88). The R squared are 0.87 (0.77), 0.92 (0.75), 0.79 (0.73),

299 0.72 (0.69), and 0.57 (0.57). It is evident that all models, except the DBN model, can fit the training set

300 data well. The DBN model exhibits the weakest performance on both the training and validation sets.

301 Alternatively, the LightGBM and XGBoost models demonstrate superior prediction performance on

302 the training set compared to the validation set. The scatter points of the training sets of these two

303 models accumulate on the 1:1 diagonal, indicating slight overfitting. The RMSE of lightGBM,

304 XGBoost, RF, DBN, and MLP models on the test set in December 2021 (January 2022) are 0.67 m/s

305 (0.64 m/s), 0.70 m/s (0.67 m/s), 0.65 m/s (0.64 m/s), 0.77 m/s (0.74 m/s), and 0.74 m/s (0.68 m/s)

306 respectively. The FA of models on the test set in December 2021 (January 2022) are 89.68 %

307 (91.11 %), 87.90 % (89.88 %), 90.64 % (91.36 %), 86.74 % (87.71 %), and 86.08 % (89.57 %). The R

308  are 0.79 (0.77), 0.77 (0.75), 0.81 (0.78), 0.71 (0.68), and 0.75 (0.74). Considering different evaluation

309  indexes, the revision effects of the five models in two months demonstrate that RMSE is that January

310  2022 is generally lower than December 2021; FA is that January 2022 is generally higher than

311  December 2021; R is that January 2022 is generally lower than December 2021. Overall, the prediction

312  performance of the five models in January 2022 surpassed that in December 2021. Furthermore, the

313  LightGBM and RF models exhibited the best performance among the five models in the two-month test

314  sets, while the DBN model had the least effective correction effect.

315      With respect to the importance of RF characteristics (Fig.6a, c), it is indisputable that the 10 m

316  wind speed predicted by WRF plays a dominant role in correcting the actual wind speed. The ones

317  following are latitude, longitude and topographic height, which represent spatial geographic

318  information, and the actual wind speed is closely related to geographic information. Subsequently,

319  relative humidity is of lesser importance. The distribution of the humidity field typically correlates with

320  the movement of the atmosphere, which is also closely related to wind speed. Certain meteorological

321  elements, such as rainfall, 2 m dew-point temperature, and 2 m temperature, contribute less importance.

322

**Table 2. The best hyper-parameters of the models**

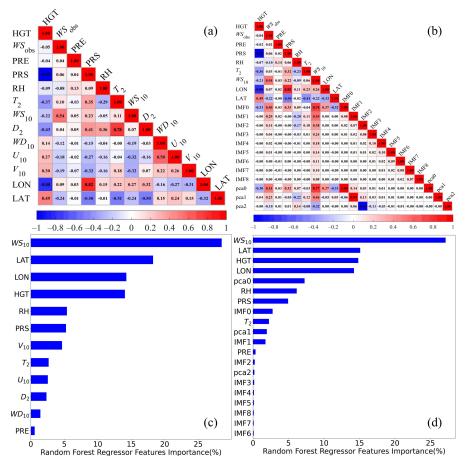| Model | parameters |
| --- | --- |
| VMD-PCA-lightGBM | 'max_depth' : 28, 'min_child_samples' : 30, 'n_estimators' : 436, 'num_leaves' : 287 |
| VMD-PCA-XGBoost | 'gamma' : 1, 'max_depth' : 19, 'min_child_weight' : 1, 'n_estimators': 408 |
| VMD-PCA-RF | 'max_depth' : 31, 'max_features' : 14, 'min_samples_leaf' : 28, 'min_samples_split' : 3, 'n_estimators' : 371 |
| VMD-PCA-DBN | 'input_length' : 20, 'output_length' : 1, 'loss_function' : 'MSE', 'optimizer' : 'Adam', 'hidden_units' : [400, 200], 'batch_size' :20000, 'epoch_pretrain' : 100, 'epoch_finetune' : 200 |
| VMD-PCA-MLP | 'batch_size' : 10114, 'hidden_layer_sizes' : 305, 'max_iter' : 386 |
| lightGBM | 'max_depth' : 21, 'min_child_samples' : 19, 'n_estimators' : 312, 'num_leaves' : 297 |

| XGBoost | 'gamma' : 0, 'max_depth' : 21, 'min_child_weight' : 9, 'n_estimators': 299 |
| RF | 'max_depth' : 40, 'max_features' : 12, 'min_samples_leaf' : 23, 'min_samples_split' : 2, 'n_estimators' : 440 |
| DBN | 'input_length' : 12, 'output_length' : 1, 'loss_function' : 'MSE', 'optimizer' : 'Adam', 'hidden_units' : [400, 200], 'batch_size' : 20000, 'epoch_pretrain' : 100, 'epoch_finetune' : 200 |
| MLP | 'batch_size' : 10232, 'hidden_layer_sizes' : 494, 'max_iter' : 311 |

323



Figure 6: Schematic diagram of correlation and feature importance for two sets of experiments. (a) and (c) represent experiment 1, and (b) and (d) represent experiment 2.

328 **3.2 Experiment 2 evaluation**

329      Experiment 2 builds upon Experiment 1, concentrating on the predicted 10-meter wind speed by

330 the WRF model. We use the VMD algorithm to decompose the predicted wind speed into 9

331 components, and use the PCA algorithm to extract the main 3 principal components. In the RF feature

332 importance analysis (Fig.6b, d), it is evident the VMD algorithm can decompose IMF0 and IMF1, with

333 contributions surpassing those of 2-meter temperature and precipitation, respectively. The importance

334 of the pca0 component, after PCA principal component extraction, reaches up to 8%. What is

335 particularly interesting is that in the correlation analysis, the correlation values between the IMF0 and

336 pca0 components and the actual wind speed are 0.50 and 0.51, which are second only to the forecasted

337 wind speed.

338      The RMSE between the predicted value and the observed value of the training set (validation set)

339 in the VMD-PCA-lightGBM, VMD-PCA-XGBoost, VMD-PCA-RF, VMD-PCA-DBN, and

340 VMD-PCA-MLP models are 0.33 m/s (0.53 m/s), 0.31 m/s (0.54 m/s), 0.52 m/s (0.57 m/s), 0.75 m/s

341 (0.75 m/s) and 0.60 m/s (0.66 m/s). The FA are 0.99 (0.94), 1.00 (0.94), 0.94 (0.93), 0.87 (0.87), and

342 0.91 (0.90). The R squared are 0.91 (0.77), 0.93 (0.77), 0.79 (0.73), 0.55 (0.55), and 0.71 (0.65). These

343 are shown in supplementary materials Figures S6-8. In comparison to the above five artificial

344 intelligence methods, training results of VMD-PCA-DBN are relatively inferior. VMD-PCA-lightGBM

345 and VMD-PCA-XGBoost models still train the processed data effectively. According to the scatter

346 density figure (Fig.7a, Fig.8a), the scatter points are relatively concentrated on the 1:1 line. The RMSE

347 of VMD-PCA-lightGBM, VMD-PCA-XGBoost, VMD-PCA-RF, VMD-PCA-DBN, and

348 VMD-PCA-MLP models on the test set in December 2021 (January 2022) are 0.63 m/s (0.63 m/s),

349 0.68 m/s (0.66 m/s), 0.62 m/s (0.64 m/s), 0.77 m/s (0.76 m/s), and 0.71 m/s (0.69 m/s) respectively.

350 The FA of the five models on the test set in December 2021 (January 2022) are 91.13 % (91.49 %),

351 89.22 % (90.23 %), 91.79 % (91.57 %), 87.93 % (87.61 %), and 87.20 % (88.94 %). The R are 0.81

352 (0.78), 0.78 (0.76), 0.82 (0.78), 0.71 (0.67), and 0.75 (0.73). The test results of the five models in

353 Experiment 2 in December 2021 and January 2022 show that the error indexes of RMSE and FA of

354 each model exhibit minimal difference in two months. Nonetheless, disregarding the correlation

355 coefficient (R) results, the performance of the five models in December 2021 is inferior to that in

356 January 2022. The diurnal variation scatter plot of two months is tested. The red scatter represents the

357    nighttime wind speed, which is more concentrated on the 1:1 line. In contrast, the blue scatter

358    represents the afternoon wind speed, which is slightly away from the 1:1 line. This suggests that the

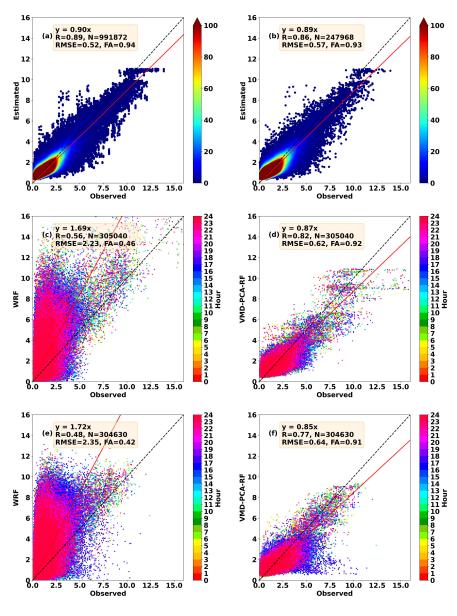359    correction effect of the five models exhibits a noticeable diurnal variation.



360

**Figure 7: The 24-hour scatter density map compared with the actual 10-meter wind speed. (a) 10-fold cross-validation training set of VMD-PCA-RF model in February 2022, (b) 10-fold cross-validation validation set of VMD-PCA-RF model in February 2022, (c) WRF forecasts in December 2021, (d)**

364     **VMD-PCA-RF model forecasts in December 2021, (e) WRF forecasts in January 2022, and (f)**

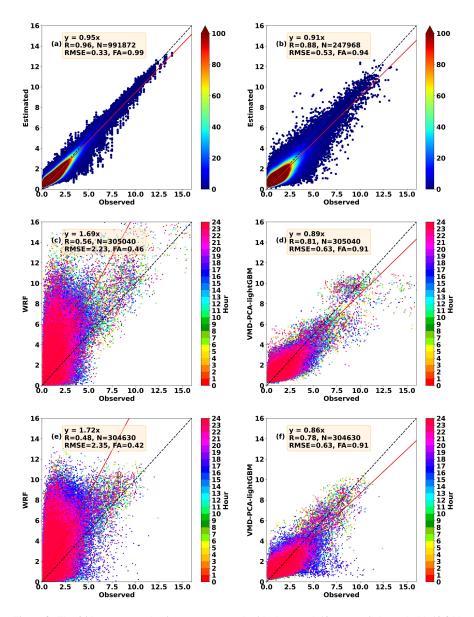365     **VMD-PCA-RF model forecasts in January 2022.**

366

**Figure 8: The 24-hour scatter density map compared with the actual 10-meter wind speed. (a) 10-fold cross-validation training set of VMD-PCA-lightGBM model in February 2022, (b) 10-fold cross-validation validation set of VMD-PCA- lightGBM model in February 2022, (c) WRF forecasts in December 2021, (d) VMD-PCA- lightGBM model forecasts in December 2021, (e) WRF forecasts in January 2022, (f) VMD-PCA- lightGBM model forecasts in January 2022.**

374 **3.3 Comparison of the two experiments**

375       Firstly, all 10 models effectively corrected the 10-meter wind speed forecasted by WRF. Table 3

376 and Table 4 represent the evaluation indexes of wind speed errors predicted by 10 models in December

377 2021 and January 2022. From the two tables, it is evident that the VMD-PCA-RF and

378 VMD-PCA-lightGBM models have the best performance in December 2021 and January 2022,

379 respectively, with the most comprehensive performance of the forecast indicators. The MAE, RMSE,

380 rMAE, rMAE, and FA for the two models VMD-PCA-RF (VMD-PCA-lightGBM) were 0.46 m/s (0.45

381 m/s), 0.62 m/s (0.63 m/s), 37.36 % (34.75 %), 50.39 % (48.65 %), and 91.79 % (91.49 %) in December

382 2021 (January 2022). Additionally, based on the analysis of the Taylor chart (Fig.9e, f) of 10 models in

383 Fig.9, it can also be seen that the scatter distance of VMD-PCA-RF and VMD-PCA-lightGBM models

384 is closest to the observed black dotted line and the black triangle position. The two models show that

385 the standard deviation is close to the observed wind speed, with the lowest RMSE and the highest R.

386 Secondly, in the comparison of cumulative probability distributions, all models passed Kolmogorov's

387 5 % confidence interval test when the interval of wind speed is 0.5 m/s (Fig.9a, d). However, when the

388 interval of wind speed is 0.2 m/s (Fig.9b, e), VMD-PCA-lightGBM model deviated from

389 Kolmogorov's 5 % confidence interval detection in December 2021. This indicates that the

390 VMD-PCA-RF model has a better predictive effect than VMD-PCA-lightGBM model in December

391 2021 when the actual wind speed is within the range of 0.4 m/s-0.8 m/s.

392     **Table 3. Table of evaluation indexes of wind speed error predicted by 10 models in December 2021**

| Model | MAE（m/s） | RMSE（m/s） | rMAE（%） | rRMSE（%） | FA（%） | R |
|---|---|---|---|---|---|---|
| VMD-PCA-lightGBM | 0.47 | 0.63 | 37.67 | 51.25 | 91.13 | 0.81 |
| VMD-PCA-XGBoost | 0.49 | 0.68 | 39.84 | 54.82 | 89.22 | 0.78 |
| VMD-PCA-RF | 0.46 | 0.62 | 37.36 | 50.39 | 91.79 | 0.82 |
| VMD-PCA-DBN | 0.53 | 0.75 | 43.32 | 61.13 | 87.93 | 0.71 |
| VMD-PCA-MLP | 0.53 | 0.72 | 43.04 | 58.47 | 87.2 | 0.75 |
| lightGBM | 0.49 | 0.67 | 39.59 | 54.16 | 89.68 | 0.79 |
| XGBoost | 0.51 | 0.70 | 41.51 | 56.64 | 87.9 | 0.77 |
| RF | 0.48 | 0.65 | 38.80 | 52.32 | 90.64 | 0.81 |
| DBN | 0.56 | 0.77 | 45.25 | 62.46 | 86.74 | 0.71 |

| | MLP | 0.55 | 0.74 | 44.65 | 60.1 | 86.08 | 0.75 |

393

394 **Table 4. Table of evaluation indexes of wind speed error predicted by 10 models in January 2022**

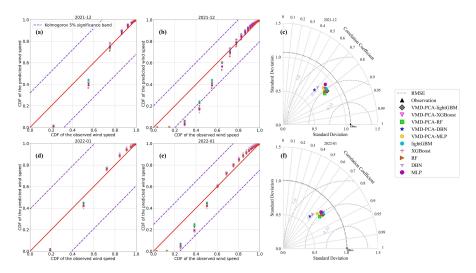| Model | MAE（m/s） | RMSE（m/s） | rMAE（%） | rRMSE（%） | FA（%） | R |
|---|---|---|---|---|---|---|
| VMD-PCA-lightGBM | 0.45 | 0.63 | 34.75 | 48.65 | 91.49 | 0.78 |
| VMD-PCA-XGBoost | 0.47 | 0.66 | 36.31 | 51.01 | 90.23 | 0.76 |
| VMD-PCA-RF | 0.46 | 0.64 | 35.06 | 49.00 | 91.57 | 0.78 |
| VMD-PCA-DBN | 0.53 | 0.75 | 40.96 | 57.49 | 87.61 | 0.67 |
| VMD-PCA-MLP | 0.50 | 0.69 | 38.46 | 53.16 | 88.94 | 0.73 |
| lightGBM | 0.46 | 0.64 | 35.24 | 49.34 | 91.11 | 0.77 |
| XGBoost | 0.48 | 0.67 | 36.68 | 51.38 | 89.88 | 0.75 |
| RF | 0.46 | 0.64 | 35.18 | 49.13 | 91.36 | 0.78 |
| DBN | 0.53 | 0.74 | 40.97 | 56.86 | 87.71 | 0.68 |
| MLP | 0.49 | 0.68 | 37.83 | 52.26 | 89.57 | 0.74 |

395



396

397 **Figure 9: The cumulative distribution probability scatter plots of the actual wind speed and the predicted**
398 **wind speed of 10 models in wind speed intervals of 0.5 m/s ((a) represents December 2021, (d) represents**
399 **January 2022) and 0.2 m/s ((b) represents December 2021, (e) represents January 2022) respectively; Taylor**
400 **distribution map ((c) represents December 2021, (f) represents January 2022).**

401

### 3.4 Spatial–temporal variations in the best models

402

403      Based on our comparative analysis results, we conclude that the best performing combination

404 models in December 2021 and January 2022 are VMD-PCA-RF and VMD-PCA-lightGBM

405 respectively. Figure 10 shows the diurnal variation corrections of the two best models for a given

406 month, as well as the diurnal variation of wind speed in the original WRF forecast. The wind speed of

407 the original WRF numerical weather forecast shows noticeable overestimation, which is confirmed in

408 Fig.8c and 8e. The scatter points of WRF forecast predominantly deviate towards the upper left corner,

409 with relatively low correlation coefficients, 0.56 and 0.23, respectively. Furthermore, the wind speed

410 forecast by WRF displays obvious diurnal variation traits, characterized by large errors between

411 afternoon and evening, specifically between 11:00 and 20:00 (Fig.10a, b). Moreover, the actual average

412 wind speed in January 2022 deviates from the range of one standard deviation of the WRF forecast

413 wind speed at 17:00 and 18:00. This demonstrates that the wind speed forecast by WRF is inaccurate

414 and exhibits substantial diurnal variation errors.

415      After the best model was corrected, the error of diurnal variation is significantly reduced (Fig.10c,

416 d). First, the average wind speed corrected by the best model is essentially consistent with the actual

417 average wind speed curve, with minimal error and no diurnal variation. Second, the one standard

418 deviation range of the corrected and actual wind speeds is also well-matched, indicating that the

419 corrected and actual wind speed distributions are consistent. The correction effect at 16:00 and 17:00

420 on January 2022 is suboptimal, which may be due to the insufficient generalization of the training

421 model and the excessive fluctuation of the actual wind speed at these two time points.

422      The FA (Fig.11a, b) and RMSE (Fig.11e, f) distribution of WRF forecast 10-meter wind speed at

423 410 stations in the five southern provinces shows that the 10-meter wind speed prediction effect of the

424 WRF model in Yunnan is superior to that in the other four provinces. In the Yunnan area, the FA of

425 most WRF forecast station 10-meter wind speeds exceeds 40 %, and RMSE value is mostly below 2.4

426 m/s. Conversely, in other regions, such as Guangxi, Guangdong and Hainan, the terrain is relatively flat.

427 The FA of the 10-meter wind speed forecast by WRF is as low as 30 % at some stations, and the

428 RMSE reaches up to 5.4 m/s. However, after the VMD-PCA-RF and VMD-PCA-lightGBM models are

429 corrected, the FA of most stations in the five southern provinces is as high as 90 %, and the RMSE is as

430 low as 0.6 m/s. Moreover, in Guangxi, Guangdong, and Hainan, where the WRF forecast effect is

431    subpar, the accuracy of the corrected 10-meter wind speed by VMD-PCA-RF (VMD-PCA-lightGBM)
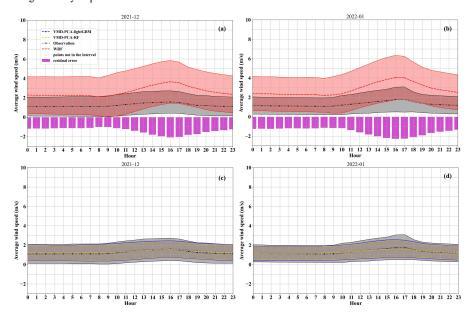
432    is significantly improved.



434    **Figure 10: VMD-PCA-lightGBM,VMD-PCA-RF and WRF daily variation of predicted and actual wind**
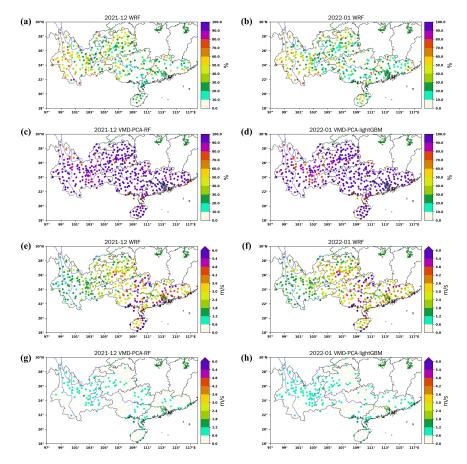435    **speeds in December 2021 and January 2022.**

**Figure 11: FA and RMSE distribution maps of VMD-PCA-RF, VMD-PCA-lightGBM and WRF models on 410 sites in five southern provinces ((a), (c), (e), and (g) represent December 2021; (b), (d), (f), and (h) represent January 2022).**

## 4. Discussion

### 4.1 The effects of BOA-VMD-PCA

It is shown in Table 2 that the hyper-parameters of the 10 models in the two experiments are different. Since the DBN model is not added to the scikit-learn Python learning package, it is challenging to call the BOA algorithm for tuning parameters. Apart from the DBN model, all the other models are optimized using the BOA algorithm. From the various evaluation indicators in Table 3 and Table 4, the DBN model, which does not use the BOA algorithm to adjust the model parameters to

25

448    obtain an optimal parameter configuration, yields the worst prediction results in December 2021 and

449    January 2022. Moreover, studies (Xiong et al., 2022) also have shown that BOA can further improve

450    the model's prediction accuracy by configuring optimal hyper-parameters. The hyper-parameters such

451    as the number of neurons and learning rate in the hidden layer, significantly impact the model's

452    performance. When the same model is applied to different data sets of two experiments, the BOA

453    adaptively obtains the optimal combination of hyper-parameters, overcoming the limitations of manual

454    parameter adjustment (Guo et al., 2021). This suggests that the selection of model hyper-parameters

455    introduces considerable uncertainty into our prediction results. Therefore, the choice of optimization

456    model parameters represents one source of uncertainty in the correction results, which entails the

457    complexity of parameter selection. However, a more advanced parameter tuning method, such as the

458    BOA tuning algorithm, is essential.

459        The VMD is used to obtain unknown but meaningful features hidden in the 10-meter wind speed

460    sequences predicted using WRF models (Li et al., 2022). In addition, the PCA can extract important

461    components of anemometer subsequences. When the stationary subsequence serves as an input to the

462    error correction model, it contains more valuable information than the previous non-stationary wind

463    speed sequences (Xu et al., 2021).

464        The complexity of the input factors in this study is one of the sources of uncertainty in the process

465    of correcting WRF prediction results. The input factors of the two experiments are not identical. In the

466    second set of experiments, the input of meteorological factors is reduced based on the first set of

467    experiments, while component information of the 10-meter wind speed predicted by WRF is increased.

468    Multiple wind speed components processed by VMD-PCA and noise reduction are introduced. Among

469    them, the importance of pca0 and IMF0 introduced is approximately 5 %. In the 10-month test sets, the

470    correction accuracy of experiment 2 is no less than the results of experiment 1 (Fig.14, Fig.S9, 10),

471    indicating that the 10-meter wind speed components introduced by the VMD-PCA contribute

472    positively to the correction results.

473

**4.2 RF feature importance**

475        In order to further understand the feature importance ranking of the RF models, we divided the

476    model prediction results and actual wind speeds of the 410 stations into 20 equal parts according to

477    height (Fig.12). First of all, the actual wind speed in December 2021 and January 2022 varies with the

478    height of the station, showing that the lower the height of the station, the more significant the change of

479    wind speed. This relationship is associated with the wind speed profile of the atmosphere, where wind

480    speed increases as height decreases. Secondly, the wind speed during the day is generally greater than

481    the wind speed at night, which is related to the turbulent motion of the atmosphere during the day.

482    Solar radiation causes the atmosphere to mix, resulting in convective movement. The 10-meter wind

483    speed at night is affected by the cooling radiation of the surface, and the atmosphere is relatively stable.

484        The 10-meter wind speed predicted by WRF has the highest feature importance in the correction

485    process of the RF models. Input factors with distinct geographic information, such as latitude,

486    longitude, and height, rank highly in feature importance. Similarly, when Sun et al. 2019 used machine

487    learning to correct the 10-meter wind speed predicted by the numerical weather prediction model

488    ECMWF, the characteristic weight of the 10-meter wind speed predicted by the model is the highest,

489    followed by the sea-land factor. Also, as the 10-meter wind speed forecast by WRF increases, the

490    instability of the 10-meter wind speed corrected by the 10 machine learning models gradually increased,

491    and the correction accuracy gradually decreased (Fig.13). This partly explains the higher importance of

492    the 10-meter wind speed forecast by WRF.

493        With 1 km as the center, the measured 10-meter wind speed is more unstable in areas where the

494    station height increases or decreases. However, the 10-meter wind speed predicted by WRF being more

495    unstable with the station height decreases (Fig.12). The VMD-PCA-RF and VMD-PCA-lightGBM

496    models significantly reduce the instability of the 10-meter wind speed predicted by WRF. When the

497    height of the station increases or decreases at 1 km, the correction intensity tends to increase gradually.

498    This further explains the higher importance of the height factor in the RF model training.
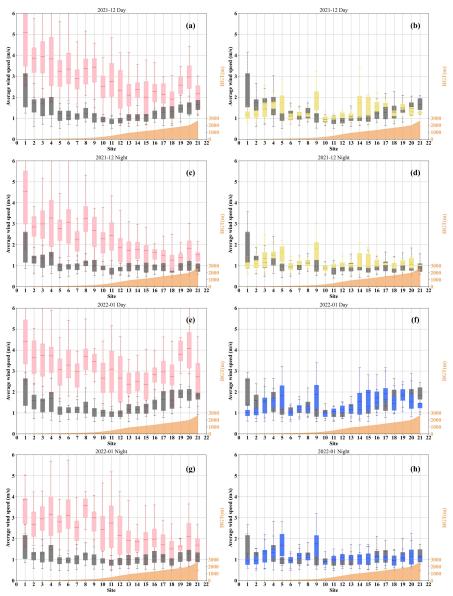
Figure 12: The boxplots of the predicted wind speeds of the VMD-PCA-RF (yellow), VMD-PCA-lightGBM (blue), and WRF (pink) models at 20 stations at different height intervals, and the boxplots of the actual wind speeds (gray).
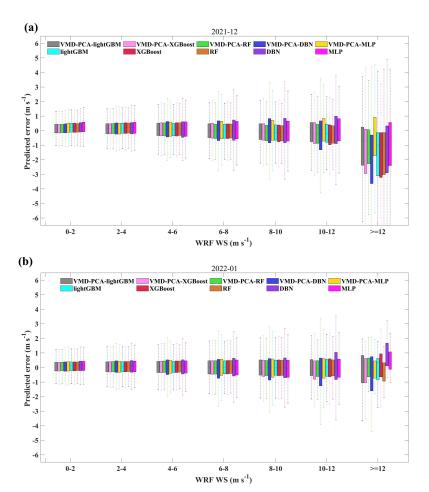
504

**Figure 13: The prediction error boxplots of 10 models in different WRF prediction intervals.**

506

### 4.3 Stability analysis of the proposed models

508    In order to identify the best model of the five southern provinces and assess the model's stability,

509    we evaluated all 10 models over 10 different months. Fig.14 shows the evaluation histogram of the

510    10-meter wind speed predicted by the 10 models in Experiment 1 and Experiment 2, as well as the

511    actual wind speed in various months. Meanwhile, Fig.S9 and S10 can more effectively illustrate the

512    daily changes of the revised results of 10 models in 10 different months. As shown in the figure 14, the

513    evaluation indexes of the model trained in Experiment 2, after VMD-PCA processing, outperform

514    those of the model trained in Experiment 1. The RF model demonstrates exceptional robustness, while

515    the MLP model exhibits the poorest performance. VMD-PCA-RF evaluation indexes are relatively

516    stable across the 10 months, with a correlation coefficient R above 0.6, accuracy rate FA above 85 %,

517    MAE below 0.6 m/s, RMSE below 0.8 m/s, rMAE below 60 %, and rRMSE below 75 %. However, the

518    robustness of the VMD-PCA-lightGBM and VMD-PCA-XGBoost models is inferior to that of the

519    VMD-PCA-RF, with all six evaluation indexes performing worse than the VMD-PCA-RF as the

520    seasons and months change. In general, VMD-PCA-RF is the best wind speed correction model for

521    winter and even throughout the entire year in the five southern provinces.
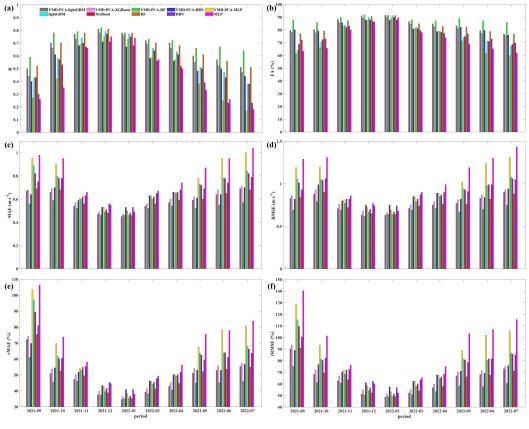


**Figure 14: Evaluation histograms of 10-meter wind speed predicted by 10 models and actual wind speed in different months in Experiment 1 and Experiment 2 ((a), (b), (c), (d), (e), and (f) represent R, FA (%), MAE (m/s), RMSE (m/s), rMAE (%), and rRMSE (%) respectively).**

527 **5. Conclusions**

528      In an effort to enhance the wind speed prediction performance for wind farms, this study

529 developed a WRF-based multi-step wind speed prediction model. A hybrid error correction strategy

530 combining BOA, VMD, PCA, and RF (LightGBM) is proposed to increase the accuracy of WRF

531 simulations. The first group of experiments used various meteorological elements as input factors in a

532 control experiment. In the second group of experiments, the wind speed sequence predicted by the

533 WRF model was decomposed into multiple IMFs using the VMD algorithm for feature extraction. A

534 principal component analysis method is used to extract meaningful principal components from these

535 subsequence IMFs to improve computational efficiency. In the error correction model, RF (lightGBM)

536 and other algorithms are used to train the relationship between different input factors and the actual

537 wind speed error, respectively.

538      Through a case analysis of 410 stations in five southern provinces in China, the following

539 conclusions can be drawn: (1) The machine learning models tuned by the BOA-VMD-PCA algorithm

540 exhibit a positive impact on wind speed error correction; (2) Feature importance analysis revealed that

541 the top eight contributing factors for correcting WRF forecasted wind speed include WRF forecast

542 10-meter wind speed (WS10), latitude, longitude, altitude, pca0, humidity, pressure, IMF0; (3)

543 VMD-PCA-RF and VMD-PCA-lightGBM are the most suitable wind speed correction algorithms for

544 December 2021 and January 2022, respectively. The MAE, RMSE, FA, rMAE, rRMSE, and R of the

545 corrected wind speed and the actual wind speed are 0.46 (0.45), 0.62 m/s (0.63 m/s), 37.36 %

546 (34.75 %), 50.39 % (48.65 %), 91.79 % (91.49 %), and 0.82 (0.78); and (4) The proposed wind speed

547 correction model (VMD-PCA-RF) demonstrates the highest prediction accuracy and stability in the

548 five southern provinces in nearly a year and at different heights. VMD-PCA-RF evaluation indexes for

549 10 months remain relatively stable: correlation coefficient R is above 0.6, accuracy rate FA is above

550 85 %, MAE is below 0.6 m/s, RMSE is below 0.8 m/s, rMAE is below 60 %, and rRMSE is below

551 75 %. In future research, the proposed VMD-PCA-RF algorithm can be extrapolated to the 3 km grid

552 points of the five southern provinces to generate a 3km grid-corrected wind speed product.

553

554

**Code availability**

556 The code and model are available as a free-access repository on Zenodo at

557 https://doi.org/10.5281/zenodo.7940686 (Zhou, 2023).

**Data Availability**

559 The data is available as a free-access repository on Zenodo at https://doi.org/10.5281/zenodo.7940686

560 (Zhou, 2023).

**Author contributions**

562 SZ developed the software, visualized the data, and prepared the original draft. SZ and YG developed

563 the methodology and carried out the formal analysis. XX and SZ validated data. SZ, YG, XX, ZD, and

564 YL reviewed and edited the text. All authors have read and agreed to the published version of the

565 paper.

**Competing interests**

567 The authors declare that they have no conflict of interest.

**Financial support**

571

572

573 **References**

574 Barthelmie, R. J., Palutikof, J. P., and Davies, T. D.: Estimation of sector roughness lengths and the

575 effect on prediction of the vertical wind speed profile, Boundary-Layer Meteorol, 66, 19–47,

576 https://doi.org/10.1007/BF00705458, 1993.

577 Cassola, F. and Burlando, M.: Wind speed and wind energy forecast through Kalman filtering of

578 Numerical Weather Prediction model output, Applied Energy, 99, 154–166,

579 https://doi.org/10.1016/j.apenergy.2012.03.054, 2012.

580 Chen, K. and Yu, J.: Short-term wind speed prediction using an unscented Kalman filter based

581 state-space support vector regression approach, Applied Energy, 113, 690–705,

582 https://doi.org/10.1016/j.apenergy.2013.08.025, 2014.

583 Cheng, W. Y. Y., Liu, Y., Liu, Y., Zhang, Y., Mahoney, W. P., and Warner, T. T.: The impact of

584 model physics on numerical wind forecasts, Renewable Energy, 55, 347–356,

585 https://doi.org/10.1016/j.renene.2012.12.041, 2013.

586 Deng, Y., Wang, B., and Lu, Z.: A hybrid model based on data preprocessing strategy and error

587 correction system for wind speed forecasting, Energy Conversion and Management, 212, 112779,

588 https://doi.org/10.1016/j.enconman.2020.112779, 2020.

589 Dhiman, H. S. and Deb, D.: A Review of Wind Speed and Wind Power Forecasting Techniques,

590 arXiv:2009.02279 [cs, eess], 2020.

591 Dong, L., Ren, L., Gao, S., Gao, Y., and Liao, X.: Studies on wind farms ultra-short term NWP wind

592 speed correction methods, in: 2013 25th Chinese Control and Decision Conference (CCDC), 2013 25th

593 Chinese Control and Decision Conference (CCDC), Guiyang, China, 1576–1579,

594 https://doi.org/10.1109/CCDC.2013.6561180, 2013.

595 Erdem, E. and Shi, J.: ARMA based approaches for forecasting the tuple of wind speed and direction,

596 Applied Energy, 88, 1405–1414, https://doi.org/10.1016/j.apenergy.2010.10.031, 2011.

597 Guo, X., Zhu, C., Hao, J., Zhang, S., and Zhu, L.: A hybrid method for short-term wind speed

598 forecasting based on Bayesian optimization and error correction, Journal of Renewable and Sustainable

599 Energy, 13, 036101, https://doi.org/10.1063/5.0048686, 2021.

600    Guo, Z., Zhao, W., Lu, H., and Wang, J.: Multi-step forecasting for wind speed using a modified

601    EMD-based    artificial    neural    network    model,    Renewable    Energy,    37,    241–249,

602    https://doi.org/10.1016/j.renene.2011.06.023, 2012.

603    Hanifi, S., Liu, X., Lin, Z., and Lotfian, S.: A Critical Review of Wind Power Forecasting

604    Methods—Past, Present and Future, Energies, 13, 3764, https://doi.org/10.3390/en13153764, 2020.

605    Hu, H., Wang, L., and Tao, R.: Wind speed forecasting based on variational mode decomposition and

606    improved    echo    state    network,    Renewable    Energy,    164,    729–751,

607    https://doi.org/10.1016/j.renene.2020.09.109, 2021.

608    Hu, J., Wang, J., and Zeng, G.: A hybrid forecasting approach applied to wind speed time series,

609    Renewable Energy, 60, 185–194, https://doi.org/10.1016/j.renene.2013.05.012, 2013.

610    Huang, Y., Yang, L., Liu, S., and Wang, G.: Multi-Step Wind Speed Forecasting Based On Ensemble

611    Empirical Mode Decomposition, Long Short Term Memory Network and Error Correction Strategy,

612    Energies, 12, 1822, https://doi.org/10.3390/en12101822, 2019.

613    Isham, M. F., Leong, M. S., Lim, M. H., and Ahmad, Z. A.: Variational mode decomposition: mode

614    determination    method    for    rotating    machinery    diagnosis,    J    VIBROENG,    20,    2604–2621,

615    https://doi.org/10.21595/jve.2018.19479, 2018.

616    James, E. P., Benjamin, S. G., and Marquis, M.: Offshore wind speed estimates from a high-resolution

617    rapidly updating numerical weather prediction model forecast dataset, Wind Energy, 21, 264–284,

618    https://doi.org/10.1002/we.2161, 2018.

619    Jiménez, P. A. and Dudhia, J.: Improving the Representation of Resolved and Unresolved Topographic

620    Effects on Surface Wind in the WRF Model, Journal of Applied Meteorology and Climatology, 51,

621    300–316, https://doi.org/10.1175/JAMC-D-11-084.1, 2012.

622    Joyce, L. and Feng Z.: Global Wind Report 2023, Global Wind Energy Council,
623    https://gwec.net/globalwindreport2023 (last access: 9 May 2023), 2023.

624    Landberg, L.: Short-term prediction of the power production from wind farms, Journal of Wind

625    Engineering    and    Industrial    Aerodynamics,    80,    207–220,

626    https://doi.org/10.1016/S0167-6105(98)00192-5, 1999.

627    Li, G. and Shi, J.: Application of Bayesian model averaging in modeling long-term wind speed

628    distributions, Renewable Energy, 35, 1192–1202, https://doi.org/10.1016/j.renene.2009.09.003, 2010.

629     Li, Y., Tang, F., Gao, X., Zhang, T., Qi, J., Xie, J., Li, X., and Guo, Y.: Numerical Weather Prediction

630     Correction Strategy for Short-Term Wind Power Forecasting Based on Bidirectional Gated Recurrent

631     Unit and XGBoost, Front. Energy Res., 9, 836144, https://doi.org/10.3389/fenrg.2021.836144, 2022.

632     Liu, H., Mi, X., and Li, Y.: An experimental investigation of three new hybrid wind speed forecasting

633     models using multi-decomposing strategy and ELM algorithm, Renewable Energy, 123, 694–705,

634     https://doi.org/10.1016/j.renene.2018.02.092, 2018.

635     Liu, Y., Wang, Y., Li, L., Han, S., and Infield, D.: Numerical weather prediction wind correction

636     methods and its impact on computational fluid dynamics based wind power forecasting, Journal of

637     Renewable and Sustainable Energy, 8, 033302, https://doi.org/10.1063/1.4950972, 2016.

638     Ma, Z., Chen, H., Wang, J., Yang, X., Yan, R., Jia, J., and Xu, W.: Application of hybrid model based

639     on double decomposition, error correction and deep learning in short-term wind speed prediction,

640     Energy Conversion and Management, 205, 112345, https://doi.org/10.1016/j.enconman.2019.112345,

641     2020.

642     Prósper, M. A., Otero-Casal, C., Fernández, F. C., and Miguez-Macho, G.: Wind power forecasting for

643     a real onshore wind farm on complex terrain using WRF high resolution simulations, Renewable

644     Energy, 135, 674–686, https://doi.org/10.1016/j.renene.2018.12.047, 2019.

645     Salcedo-Sanz, S., Ángel M. Pérez-Bellido, Ortiz-García, E. G., Portilla-Figueras, A., Prieto, L., and

646     Paredes, D.: Hybridizing the fifth generation mesoscale model with artificial neural networks for

647     short-term     wind     speed     prediction,     Renewable     Energy,     34,     1451–1457,

648     https://doi.org/10.1016/j.renene.2008.10.017, 2009.

649     Salcedo-Sanz, S., Ortiz-García, E., Pérez-Bellido, Á., Portilla-Figueras, A., and Prieto, L.: Short term

650     wind speed prediction based on evolutionary support vector regression algorithms, Expert Syst. Appl.,

651     38, 4052–4057, https://doi.org/10.1016/j.eswa.2010.09.067, 2011.

652     Sun, Q., Jiao, R., Xia, J., Yan, Z., Li, H., Sun, J., Wang, L., and Liang, Z.: Adjusting Wind Speed

653     Prediction of Numerical Weather Forecast Model Based on Machine Learning Methods.

654     Meteorological Monthly, 45(3): 426-436. https://doi.org/10.7519/j.issn.1000-0526.2019.03.012, 2019.

655     Tang, R., Ning, Y., Li, C., Feng, W., Chen, Y., and Xie, X.: Numerical Forecast Correction of

656     Temperature and Wind Using a Single-Station Single-Time Spatial LightGBM Method, Sensors, 22,

657     193, https://doi.org/10.3390/s22010193, 2021.

658    Tascikaraoglu, A. and Uzunoglu, M.: A review of combined approaches for prediction of short-term

659    wind speed and power, Renewable and Sustainable Energy Reviews, 34, 243–254,

660    https://doi.org/10.1016/j.rser.2014.03.033, 2014.

661    Wang, C., Zhang, H., Fan, W., and Ma, P.: A new chaotic time series hybrid prediction method of wind

662    power based on EEMD-SE and full-parameters continued fraction, Energy, 138, 977–990,

663    https://doi.org/10.1016/j.energy.2017.07.112, 2017.

664    Wang, J. and Hu, J.: A robust combination approach for short-term wind speed forecasting and analysis

665    – Combination of the ARIMA (Autoregressive Integrated Moving Average), ELM (Extreme Learning

666    Machine), SVM (Support Vector Machine) and LSSVM (Least Square SVM) forecasts using a GPR

667    (Gaussian Process Regression) model, Energy, 93, 41–56, https://doi.org/10.1016/j.energy.2015.08.045,

668    2015.

669    Williams, J. L., Maxwell, R. M., and Monache, L. D.: Development and verification of a new wind

670    speed forecasting system using an ensemble Kalman filter data assimilation technique in a fully

671    coupled hydrologic and atmospheric model: Data Assimilation in a Coupled Forecasting System, J.

672    Adv. Model. Earth Syst., 5, 785–800, https://doi.org/10.1002/jame.20051, 2013.

673    Xiong, X., Guo, X., Zeng, P., Zou, R., and Wang, X.: A Short-Term Wind Power Forecast Method via

674    XGBoost    Hyper-Parameters    Optimization,    Front.    Energy    Res.,    10,    905155,

675    https://doi.org/10.3389/fenrg.2022.905155, 2022.

676    Xu, Q., He, D., Zhang, N., Kang, C., Xia, Q., Bai, J., and Huang, J.: A Short-Term Wind Power

677    Forecasting Approach With Adjustment of Numerical Weather Prediction Input by Data Mining, IEEE

678    Trans. Sustain. Energy, 6, 1283–1291, https://doi.org/10.1109/TSTE.2015.2429586, 2015.

679    Xu, W., Liu, P., Cheng, L., Zhou, Y., Xia, Q., Gong, Y., and Liu, Y.: Multi-step wind speed prediction

680    by combining a WRF simulation and an error correction strategy, Renewable Energy, 163, 772–782,

681    https://doi.org/10.1016/j.renene.2020.09.032, 2021.

682    Zhang, D., Peng, X., Pan, K., and Liu, Y.: A novel wind speed forecasting based on hybrid

683    decomposition and online sequential outlier robust extreme learning machine, Energy Conversion and

684    Management, 180, 338–357, https://doi.org/10.1016/j.enconman.2018.10.089, 2019a.

685    Zhang, Y., Chen, B., Pan, G., and Zhao, Y.: A novel hybrid model based on VMD-WT and

686    PCA-BP-RBF neural network for short-term wind speed forecasting, Energy Conversion and

687    Management, 195, 180–197, https://doi.org/10.1016/j.enconman.2019.05.005, 2019b.

688    Zhang, Z., Ye, L., Qin, H., Liu, Y., Wang, C., Yu, X., Yin, X., and Li, J.: Wind speed prediction

689    method using Shared Weight Long Short-Term Memory Network and Gaussian Process Regression,

690    Applied Energy, 247, 270–284, https://doi.org/10.1016/j.apenergy.2019.04.047, 2019c.

691    Zhao, J., Guo, Z.-H., Su, Z.-Y., Zhao, Z.-Y., Xiao, X., and Liu, F.: An improved multi-step forecasting

692    model based on WRF ensembles and creative fuzzy systems for wind speed, Applied Energy, 162,

693    808–826, https://doi.org/10.1016/j.apenergy.2015.10.145, 2016.

694    Zhao, J., Wang, J., Guo, Z., Guo, Y., Lin, W., and Lin, Y.: Multi-step wind speed forecasting based on

695    numerical simulations and an optimized stochastic ensemble method, Applied Energy, 255, 113833,

696    https://doi.org/10.1016/j.apenergy.2019.113833, 2019.

697    Zhou, S.: A hybrid method for numerical weather prediction wind speed based on Bayesian

698    optimization (version 1.2.0) and error correction: First release of my code. Zenodo [code]

699    https://doi.org/10.5281/zenodo.7940686, 2023.

700    Zjavka, L.: Wind speed forecast correction models using polynomial neural networks, Renewable

701    Energy, 83, 998–1006, https://doi.org/10.1016/j.renene.2015.04.054, 2015.

702