1

From Numerical Weather Prediction outputs to accurate local surface wind speed : statistical modeling and forecasts.

Alonzo B.^{1,2}, Plougonven R.¹, Mougeot M.², Fischer A.², Dupré A.¹ and Drobinski P.¹

¹LMD/IPSL, cole Polytechnique, Université Paris Saclay, ENS, PSL Research University, Sorbonne Universités, UPMC Univ Paris 06, CNRS, Palaiseau France

²Laboratoire de Probabilités et Modèles Aléatoires, Université Paris Diderot - Paris 7, Paris, France.

November 8, 2017

Abstract

Downscaling a meteorological quantity at a specific location from outputs of Numerical Weather Prediction models is a vast field of research with continuous improvement. The need to provide accurate forecasts of the surface wind speed at specific locations of wind farms has become critical for wind energy application. While classical statistical methods like multiple linear regression have been often used in order to reconstruct wind speed from Numerical Weather Prediction model outputs, machine learning methods, like Random Forests, are not as widespread in this field of research. In this paper, we compare the performances of two downscaling statistical methods for reconstructing and forecasting wind speed at a specific location from the European Center of Medium-range Weather Forecasts (ECMWF) model outputs. The assessment of ECMWF shows for 10m wind speed displays a systematic bias, while at 100m, the wind speed is better represented. Our study shows that both classical and machine learning methods lead to comparable results. However, the time needed to pre-process and to calibrate the models is very different in both cases. The multiple linear model associated with a wise pre-processing and variable selection shows performances that are slightly better, compared to Random Forest models. Finally, we highlight the added value of using past observed local information for forecasting the wind speed on the short term.

Keywords

Local wind speed, Downscaling, Statistical modeling, Numerical Weather Prediction model, Wind speed forecasts.

1 **Introduction**

The wind energy sector has seen a very sharp growth in recent years. Accord-2 ing to the Global Wind Energy Council (GWEC), 54GW has been installed 3 in 2016, corresponding to an increase of 12.6% of the total installed capacity 4 ([GWE, 2016]). Worldwide, the number of wind farms increases each year 5 and feeds the electrical network with a larger amount of energy. For instance, 6 in 2016, France has seen its highest capacity growth rate ever recorded. This 7 sharp increase of connected wind power has for instance allowed the network 8 to receive 8.6 GW from wind power plants, on November 20th, correspond-9 ing to 17.9% of the energy produced this day ([RTE, 2016]). The need to 10 have access to accurate wind forecasts on several timescales is thus becoming 11 crucial for the wind energy producer and grid operator, in order to antici-12 pate the energy production, to plan maintenance operations and to handle 13 balance between energy production and consumption. Changing regulations 14 of the energy market with the end of feeding-in tariffs make this anticipa-15 tion vital for wind energy producers. Finally, a related but different topic 16 consists in the estimation of the wind resource of its long-term (multi-year) 17 variability and trends mainly for prospecting purposes. 18

The increasing need for accurate forecasts of the surface wind speed fortunately comes with the improvement of the Numerical Weather Prediction models (NWP) describing and forecasting atmospheric motions. Indeed, they constitute a key source of information for surface wind speed forecasts all the more so as their realism, accuracy and resolution have increased steadily over the years [Bauer et al., 2015].

Nevertheless, these models are not necessarily performing uniformly well 25 for all atmospheric variables. Their astonishing performances are evalu-26 ated on variables such as mid-tropospheric pressure which reflect the large-27 scale mass distribution, which is effectively well understood physically (e.g. 28 [Vallis, 2006]) and efficiently modeled numerically. Variables tied to phenom-29 ena occurring on smaller scales (such as cloud-cover or near-surface winds) 30 depend much more directly on processes that are *parameterized* (e.g. not 31 resolved). In contrast to large-scale motions (governed by the Navier-Stokes 32 equations), parameterizations are generally partly rooted in physical argu-33 ments, but also in large part empirical. When comparing output from a 34 numerical model to a local measurement, there will therefore always be sev-35 eral sources of error: representativity error (contrast between the value over 36 a grid-box and the value at a specific point), numerical error (even if we 37 were describing only processes governed by well-established physical laws, 38 discretization is unavoidable), and error tied to the physics described (be-39 cause processes, especially parameterized ones, are not well modeled). To 40 reduce representativity error and to better represent small-scale processes, in 41 particular those tied to topography and surface roughness, one strategy con-42 sists in downscaling with models that describe the atmospheric flow on finer 43 scales (e.g [Wagenbrenner et al., 2016]). One disadvantage of this approach 44 is the numerical cost, and one limitation is the need for finer observations to 45 initialize the state of the atmosphere, if details of the flow other than those 46 directly implied by the topography and surface condition are sought for. 47

48 Strategies to estimate surface winds, or other meteorological variables,

from the output of Numerical Weather Prediction models (NWP) or climate
models have been developed in several contexts, with different motivations,
and leading to different methodologies and applications.

Model Output Statistics (MOS) has been developed in weather forecast-52 ing for several decades to estimate the *weather related* variability of a physi-53 cal quantity, based on NWP model output [Glahn and Lowry, 1972]. NWP 54 models perform now very well in predicting large-scale systems. Relations 55 thus can be derived to link the latters to local variables at an observation 56 site. Linear models are generally used, with the outcome now expanded 57 over a wider area than only the location of stations where observations are 58 available [Zamo et al., 2016]. 59

In the context of climate change, downscaling a meteorological quan-60 tity at a given location in order to produce time series which have plau-61 sible statistical characteristics under climate change has for long been in-62 vestigated [Wilby et al., 1998]. The challenge is here to capture appropri-63 ately the relation between large-scale flow (as it can be described by a 64 model with a moderate or low resolution) and a variable at a specific lo-65 cation (e.g. wind, temperature, precipitation) and then use climate mod-66 els to provide a description of the large-scale atmospheric state under cli-67 mate change. Local time series with appropriate variability and consis-68 tent with this large-scale state of the atmosphere are then generated, e.g. 69 [Salameh et al., 2009, Maraun et al., 2010, Wilby and Dawson, 2013]. 70

Wind energy domain is nowadays a very active branch in downscal-71 ing techniques because of the need for accurate forecasts at specific loca-72 tion of a wind farm. For describing winds close to the surface, 10m wind 73 speed is often a convenient variable as it has been for decades a reference 74 observed variable and also now a reference NWP model output. In the 75 case of wind energy, the wind speed then needs to be extrapolated at the 76 hub height to have access to wind power, leading to an increase of the er-77 ror on the predicted power ([Kubik et al., 2011], [Howard and Clark, 2007], 78 [Mohandes et al., 2011]). Wind speed at the hub height (typically 100m) 79 is a variable of interest as it allows to avoid vertical extrapolation errors 80 ([Cassola and Burlando, 2012]), but it is rarely available in observations. 81 Different outputs of NWP models can be used as explanatory variables of 82 the near surface wind speed. It seems that there is no strong consensus 83 on the predictors to use, mainly because relations between predictors and 84 predict and should differ from one location to the other. However, different 85 studies have shown the importance of a certain set of variables to predict 86 surface wind speed. Amongst them, markers of large-scale systems (geopo-87 tential height, pressure fields) and boundary layer stability drivers (surface 88 temperature, boundary layer height, wind and temperature gradient) can 89 be cited ([Salameh et al., 2009], [Devis et al., 2013], [Davy et al., 2010]). In 90 terms of methodology, several models have already been studied, including 91 Linear regression, Support Vector Models (SVM) or Artificial Neural Net-92 work (ANN) (Jung and Broadwater, 2014], [Soman et al., 2010]). 93

The model of the European Center for Medium-range Weather Forecasts (ECMWF) has reached a resolution of about 9km in the horizontal. In addition, ECMWF analyses and forecasts now give access to 100m wind speed output, developed mainly for wind energy applications. If we can be very

confident in the ability of NWP models to represent several variables, some 98 others may not be so reliable. This is especially the case for surface variables 99 such as 10m and 100m wind speed. Consequently, using the robust informa-100 tion given by some variables to correct surface wind speed is straightforward. 101 We have access to surface wind speed observed at 10m, 100m over a long 102 period of 5 years at SIRTA observation platform [Haeffelin et al., 2005]. The 103 aim of this project is, in particular, to explore how different statistical mod-104 els perform in forecasting the 10m and 100m wind speed using informations 105 of ECMWF analyses and forecasts outputs at different horizons. We choose 106 multiple linear regression because it is a widely used technique, and Ran-107 dom Forests which have not been, to our knowledge, deeply studied in the 108 framework of downscaling surface wind speed. For multiple linear regres-109 sion, variable selection is a very important step for calibrating the statistical 110 models, whereas Random Forests handle variables automatically. Moreover, 111 Random Forests can handle nonlinear relations very well. Therefore, the 112 comparison of those very different statistical models, as well as the informa-113 tion used by each of them, should be very instructive. 114

The paper is organized in 5 parts. The next section describes together the data and the statistical models to be used. In section 3, the training dataset is explored, and used to calibrate the statistical models. In section 4, forecasts of 10m and 100m wind speed are run to downscale wind speed at the observation site. In the last section, we discuss the results, conclude and give perspectives to this work.

¹²¹ 2 Data and Methodology

122 2.1 Data

123 Observed Wind speed

In this paper, we use observations of the wind speed at the SIRTA obser-124 vation platform ([Haeffelin et al., 2005]). Surface wind speed at 10m height 125 from anemometer recording is available at the 5-minutes frequency. The 126 wind speed at 100m height from Lidar recording is available at 10-minutes 127 frequency. Both data span for 5 years from 2011 to 2015. We filter obser-128 vations by a sinusoidal function over a 6-hour window centered at 00h, 06h, 129 12h and 18h to obtain a 6-hourly observed wind speed to be compared to 130 the NWP model outputs available at this time frequency. We found that 131 the resulting time series are not sensitive to the filter function. We also try 132 different filtering windows, concluding that 6-hours is the best to compare 133 to the NWP model outputs. Due to some missing data, two final time se-134 ries of 5049 filtered observations are computed (over 7304 if all data were 135 available). 136

SIRTA observatory is based 20km in the South of Paris on the Saclay plateau (48.7° N and 2.2° E). Figure 1 shows the SIRTA observation platform location, marked by the red point on the map, and its close environment. Regarding the relief near SIRTA, observe that a forest is located at about 50m north to the measurement devices. South, buildings can be found at about 300m from the SIRTA observatory. In the East-West axis, no close obstacle are encountered. Further south, the edge of the Saclay plateau



Figure 1: Map of the SIRTA observation platform and its surroundings.

shows a vertical drop of about 70m, from 160m on top to 90m at the bottom.

145 NWP model outputs - ECMWF Analyses

Variables are retrieved from ECMWF analyses at 4 points around the 146 SIRTA platform. The spatial resolution of ECMWF analyses is of about 147 $16 \text{km} (0.125^{\circ} \text{ in latitude and longitude})$. Topography is thus smoothed com-148 pared to the real one. As the surface wind speed is very influenced by the 149 terrain, the modeled surface wind speed is not necessarily close to the ob-150 served wind speed. The data spans from the 01/01/2011 to 31/12/2015 at 151 the 6-hourly frequency. It is sampled at each date where a filtered sampled 152 observation is available. 153

The near surface wind speed at a given location can be linked to different 154 phenomena. The large-scale circulation brings the flow to the given location 155 explaining the slowly varying wind speed. The wind speed in altitude, the 156 geopotential height, the vorticity, the flow divergence, sometimes the tem-157 perature can be markers of large systems like depressions, fronts, storms, or 158 high pressure systems explaining a large part of the low frequency variations 159 of the surface wind speed (Table 2). At a finer scale, what is happening in 160 the boundary layer is very important to explain the intra-day variations of 161 the wind speed. The state and stability of the boundary layer can be derived 162 from surface variables describing the exchanges inside the layer. Exchanges 163 are driven mostly by temperature gradient and wind shear that develop tur-164

5

bulent flow (Table 3). Thermodynamical variables like surface, skin, and dew 165 point temperatures and surface heat fluxes can also inform on the stability 166 of the boundary layer, as well as its height and dissipation on its state (Table 167 1). In the end, 20 output variables are retrieved from ECMWF analyses at 168 the 4 points around the SIRTA observatory and at different pressure levels. 169 Note that we restrict the study to local variables (at the location of measure-170 ments or in the column above). It might also be possible to take advantage 171 from larger scale information ([Zamo et al., 2016], [Davy et al., 2010]). 172

Altitude (m)	Variable	Unit	Name
10m/100m	Norm of the Wind speed	$m.s^{-1}$	F
10m/100m	Zonal Wind speed	$m.s^{-1}$	U
10m/100m	Meridional Wind speed	$m.s^{-1}$	V
2m	Temperature	K	Т
2m	Dew point Temperature	K	Dp
Surface	Skin temperature	K	skt
Surface	mean sea level pressure	Pa	msl
Surface	Surface pressure	Pa	$^{\mathrm{sp}}$
-	Boundary layer height	m	blh
-	Boundary layer dissipation	$J.m^{-2}$	bld
Surface	Surface latent heat flux	$J.m^{-2}$	slhf
Surface	Surface sensible heat flux	$J.m^{-2}$	sshf

Table 1: Surface Variables

Pressure level (hPa)	Variable	Unit	Name
1000hPa/925hPa/850hPa/700hPa/500hPa	Zonal Wind speed	$m.s^{-1}$	U
1000hPa/925hPa/850hPa/700hPa/500hPa	Meridional Wind speed	$m.s^{-1}$	V
1000hPa/925hPa/850hPa/700hPa/500hPa	Geopotential height	$m^2.s^{-2}$	Z
1000hPa/925hPa/850hPa/700hPa/500hPa	Divergence	s^{-1}	Di
1000hPa/925hPa/850hPa/700hPa/500hPa	Vorticity	s^{-1}	Vo
1000hPa/925hPa/850hPa/700hPa/500hPa	Temperature	K	Т

 Table 2: Altitude Variables

Pressure level (hPa)	Variable	Unit	Name
10m to 925hPa	Wind shear	$m.s^{-1}$	ΔF
10m to 925hPa	Temperature gradient	K	ΔT

Table 3:	Computed	Variables
----------	----------	-----------

173 ECMWF deterministic forecasts

The year 2015 of deterministic forecasts is retrieved from ECMWF model. A forecast is launched every day at 00:00:00 UTC. The time resolution retained is of 3 hours and the maximum lead-time is 5 days. The same variables as for the analyses are retrieved at the same points around the SIRTA platform.

179 2.2 Methodology

Our aim is to model the real observed wind speed from outputs of NWP 180 model described above. More specifically, we use ECMWF analyses i.e the 181 best estimate of the atmospheric state at a given time using a model and 182 observations ([Kalnay, 2003]). In what follows, the observed wind speed is 183 the target and the analysed variables are potential explanatory features. Be-184 cause of the complexity of meteorological phenomena, statistical modeling 185 provides an appropriate framework for corrections of representativity errors 186 and the modeling of site dependent variability. In this context, two main 187 directions may be as usual investigated, parametric and nonparametric mod-188 els. 189

Parametric models assume that the underlying relation between the target variable and the explanatory variables has, relatively to a certain noise, a particular analytical shape depending on some parameters, which need to be estimated through the data. Among this family of models, the linear model with a Gaussian noise is widely used, mostly thanks to its simplicity [Friedman et al., 2001]. Associated to an adequate variable selection, it may be very effective.

Nonparametric models do not suppose in advance a specific relation be-197 tween the variables: instead, they try to learn this complex link directly from 198 199 the data itself. As such, they are very flexible, but their performance usually strongly depends on regularization parameters. The family of nonparamet-200 ric models is quite large: among others, one may cite the nearest neighbors 201 rule, the kernel rule, neural networks, support vector machines, regression 202 trees, random forests... Regression trees, which have the advantage of be-203 ing easily interpretable, show to be particularly effective when associated 204 to a procedure allowing to reduce their variance as for the Random Forest 205 Algorithm. 206

Let us describe the linear model and random forests in our context with more details. The linear model supposes a relation between the target Y_t , observed wind speed at time t, and explanatory variables X_t^1, \ldots, X_t^d , available from the ECMWF, at this time t. For lightening the notation, we omit the index t in the next equation. The linear model is given by

$$Y = \beta_0 + \sum_{j=1}^d \beta_j X_j + \varepsilon,$$

where the β_i 's are coefficients to be estimated using least-square criterion 207 minimization method, and $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ represents the noise. Among the 208 meteorological variables X_1, \ldots, X_d , some of them provide more important 209 information linked to the target than others, and some of them may be 210 correlated. In this case, the stepwise variable selection method is useful to 211 keep only the most important uncorrelated variables [Friedman et al., 2001]. 212 Denoting by $\hat{\beta}_0, \ldots, \hat{\beta}_d$ the final coefficients obtained this way (some of them 213 are zero), the estimated wind \hat{Y} is then given by 214

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^d \hat{\beta}_j X_j.$$
(1)

An alternative approach to perform variable selection and regularization is to use the Lasso method (see for instance [Tibshirani, 1994]), relying on minimization of the least square criterion penalized by the ℓ^1 norm of the coefficients β_1, \ldots, β_d . More specifically, for this model, the predicted wind speed at time t is a linear combination of all the previous variables as in equation (1), the coefficients β_1, \ldots, β_d being estimated using the least square procedure, under the constraint $\sum_{j=1}^d |\beta_j| \leq \kappa$ for some constant $\kappa > 0$.

Regression trees are binary trees built by choosing at each step the cut minimizing the intra-node variance, over all explanatory variables X_1, \ldots, X_d and all possible thresholds (denoted by S_j hereafter). More specifically, the intra-node variance, usually called deviance, is defined by

$$D(X_j, S_j) = \sum_{X_j < S_j} (Y_s - \overline{Y}^-)^2 + \sum_{X_j \ge S_j} (Y_s - \overline{Y}^+)^2,$$

where \overline{Y}^{-} (respectively \overline{Y}^{+}) denotes the average of the observed wind speed in the area $\{X_j < S_j\}$ (respectively $\{X_j \ge S_j\}$). Then, the selected j_0 variable and associated threshold is given by $(X_{j_0}, S_{j_0}) = \arg\min_{j, S_j} D(X_j, S_j)$. The prediction is provided by the value associated to the leaf in which the observation falls.

To reduce variance and avoid over-fitting, it may be interesting to generate several bootstrap samples, fitting then a tree on every sample and averaging the predictions, which leads to the so-called Bagging procedure [Breiman, 1996]. More precisely, for *B* bootstrap samples, the predicted power is given by

$$\hat{Y} = \sum_{b=1}^{B} \hat{Y}^b,\tag{2}$$

where \hat{Y}^{b} denotes the wind speed predicted by the regression tree associated with the *b*-th bootstrap sample. To produce more diversity in the trees to be averaged, an additional random step may be introduced in the previous procedure, leading to Random Forests, where the best cut is chosen among a smaller subset of randomly chosen variables. The predicted value is the mean of the predictions of the trees, as in equation (2).

²³⁷ 3 The relationship between analysed and observed ²³⁸ winds

$_{239}$ 3.1 10m/100m wind speed variability comparison

In this section we compare the observed wind speed at 10m and 100m with 240 the 10m and 100m wind speed output of the ECMWF analyses, respectively. 241 Figure 2 shows the Probability Density Function (PDF) of the wind speed 242 coming from ECMWF analyses and observations, and also for illustration 243 an example of a time series of corresponding wind speeds. It appears that 244 the 10m wind speed from ECMWF analyses displays a systematic bias by 245 overestimating the 10m observed wind speed (Figure 2, a and b). The wind 246 at 100m is comparatively well modeled in terms of variations in the time 247 series, but also in terms of distribution (Figure 2, c and d). It seems that the 248

9

errors mainly come from the overestimation of peaked wind speeds and the underestimation of low wind speeds (Figure 2, c and d). As 10m wind speed is very influenced by even low topography and surrounding obstacles, which are smoothed or not represented in ECMWF analyses, some of its variations are not well described, and even a quite systematic bias is displayed. The effect of the topography and terrain specificity have less impact on the 100m wind speed, so that it is much better represented in ECMWF analyses.



Figure 2: 10m (top) and 100m (bottom) wind speed time series in summer 2011 (panels a and c, respectively) and the respective probability density function estimated over the 5 years sample wind speed (panels b and d).

Periods	Deviation		RMSE		Correlation	
	F10	F100	F10	F100	F10	F100
2011-2015	-1.00	0.14	1.41	1.01	0.82	0.93
2011	-1.19	0.04	1.59	1.06	0.80	0.91
2012	-0.94	0.23	1.31	1.03	0.85	0.92
2013	-1.13	0.06	1.52	0.93	0.82	0.94
2014	-0.88	0.26	1.30	1.00	0.80	0.93
2015	-0.87	0.14	1.30	0.97	0.82	0.94
Winter	-0.97	0.04	1.41	0.97	0.83	0.94
Spring	-1.11	0.27	1.56	1.02	0.71	0.90
Summer	-0.92	0.33	1.31	1.04	0.80	0.91
Fall	-1.04	-0.10	1.36	1.00	0.87	0.93

Table 4: Deviation, RMSE, and Correlation performed by ECMWF analyses for modeling the 10m and 100m wind speed.

The ability of the model to represent the observed wind speed is quantified in Table 4 by the deviation, the Root Mean Square Error (RMSE), and

10

the Pearson correlation which formula are given by equations (3),(4), and (5) respectively.

Deviation for the
$$i^{th}$$
 observation = $(y_i - x_i)$ (3)

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^{n} (x_i - y_i)^2}{n}}$$
(4)

Correlation =
$$\frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n} (y_i - \bar{y})^2}},$$
 (5)

where x_i is the wind speed from the NWP model and y_i the observed wind speed ; n is the number of samples (x_i, y_i) and $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ (the sample mean) and analogously for \bar{y} .

263

No clear improvement of the ECMWF analyses over the years from 2011 264 to 2015 can be detected in Table 4. The correlation stays quite constant over 265 the years for both 10m and 100m wind speeds. The Deviation and RMSE 266 seem to decrease for the 10m wind speed but it cannot be confirmed because 267 of the good score performed in 2012. The variations of performance may 268 only come from changes in the predictability of the weather over Europe 269 [Folland et al., 2012]. Seasonal variations of the performance of ECMWF 270 analyses can be seen, especially on the correlation between the observed and 271 modeled wind speed. At both 10m and 100m, the analysed wind speed is 272 better correlated with the observations in winter and fall than in spring and 273 summer. In all cases, the scores shown are better for the 100m wind speed 274 than for the 10m wind speed. 275



Figure 3: 10m wind speed from ECMWF analyses function of the 10m observed wind speed given cardinal directions. Panels correspond to a direction modeled by ECMWF analyses; the wind blows from a. West, b. Southwest, c. South, d. Southeast, e. East, f. Northeast, g. North, h. Northwest.

Variations of the performance of the ECMWF analyses in representing the observed wind speed are evidenced by Figure 3. The figure shows the

10m wind speed from ECMWF analyses as a function of the 10m observed 278 wind speed for different directions of the analysed wind speed. It is obvious 279 that the errors made by the numerical model differ regarding the direction 280 of the wind. For instance, when the wind comes from the West (figure 3, 281 a), the correlation is very close to one, but for a wind coming from the 282 North/Northeast (Figure 3, f and e), it is very low, and the model highly 283 overestimates the 10m wind speed. It can be easily linked to the specificity of 284 the terrain. Indeed, when a northerly wind is recorded, it has been blocked 285 by the forest north of the anemometer. The same happened for southerlies 286 with the building situated further and which influence is thus not as substan-287 tial as the forest. Figure 4 displays the same as Figure 3 but for the 100m 288 wind speed. It seems that there is no more dependence of the performance 289 of the ECMWF analyses regarding the direction of the 100m wind speed; 290 it appears to be not significantly impacted by the surrounding forests and 291 buildings. 292



Figure 4: Same as Figure 3 but for 100m wind speed.

3.2 Reconstruction of the 10m/100m observed wind speed using NWP outputs

In the sequel, a k-fold cross validation is performed over 10 different periods 295 taken within the 5-years of analyses and observation. Each time, statisti-296 cal downscaling models are trained on a given period and applied over the 297 remaining period to reconstruct the 10m and 100m wind speed, so that it 298 results in 10 reconstructions that span the 5 years of data. Table 5 enu-299 merates the statistical downscaling models assessed in this study. Models 300 differ by their types (Linear Regression and Random Forests), the explana-301 tory variable selection, and whether a model is conditionally fitted regarding 302 the direction of the wind speed or not. We evaluate the different statistical 303 models in terms of RMSE and Correlation with the observed wind speed on 304 the reconstruction period. 305

306 10m wind speed reconstruction

12 FOREWER2017, 011, v1: 'From Numerical Weather Prediction outputs to accurate local...

From model output to surface wind

Model type	Explanatory variables	Direction dependance	Name
Linear	F10	No	LR_F
Linear	All	No	LR_A
Linear	Stepwise	No	LR_{SW}
Linear	Lasso	No	LR_{La}
Linear	F10	Yes	LR_F^{dir}
Linear	All	Yes	LR_A^{dir}
Linear	Stepwise	Yes	LR_{SW}^{dir}
Random Forest	All	No	RF_A
Random Forest	All	Yes	RF_{Λ}^{dir}

Table 5: Statistical models used to downscale 10m and 100m wind speed.



Figure 5: RMSE and Correlation results when reconstructing 10m wind speed with models described in Table 5. The first boxes stand for the ECMWF analyses 10m wind speed.

Figure 5 shows results for the reconstruction of the 10m wind speed. Each box contains the 10th reconstructed k-fold periods. First, by using only wind speed with a linear model LR_F , RMSE is reduced by about 40%,

13

but the correlation stays constant. Adding other variables to linear model 310 (i.e. LR_A , LR_{SW} and LR_{La}) allows to reduce the RMSE by 60%, and to sig-311 nificantly improve correlation from 0.80 to 0.91 between reconstructed wind 312 speed and observed one. Using stepwise selection of variables, the Lasso pe-313 nalization or all variables does not change results in this case, showing that 314 only a part of the information is useful. Using variable selection as stepwise 315 or ℓ_1 penalization (Lasso) avoids over-fitting. Random Forests models per-316 form slightly better than linear models without defining one given model per 317 cardinal wind directions. Variables selected stepwise are very diverse (wind 318 speed, large scale variables, boundary layer state drivers), while Lasso tech-319 nique mainly selects wind speed and wind component, thus using redundant 320 information. Analyzing the main variables used by Random Forests shows 321 that this methods seems to put much weight on wind component first, high-322 lighting the dependence of the error on the 10m wind speed regarding its 323 direction. 324



Figure 6: Timeseries (left) and PDF (right) of the observed 10m wind speed (straight black line), and 10m wind speed from ECMWF (dotted black line) (a and b), Linear models (LR_{SW} (blue) and LR_{SW}^{dir} (red)) (c and d), Random Forest models (RF_A (cyan) and RF_A^{dir} (magenta)) (e and f).

By fitting a linear model in each direction (noted with 'dir') we manually 325 introduce a relevant information, especially for 10m wind speed (Figure 3). 326 The model is however more exposed to under-fitting as the sample size of the 327 training data in one direction can be low. Nevertheless, LR_{SW}^{dir} performs bet-328 ter than all other models. Indeed, stepwise choice is made for each direction 329 so that the model is deeply adapted to each direction. This method results 330 in a significant improvement of the RMSE and correlation scores. Fitting a 331 Random Forest in each direction does not improve results, maybe because 332 the direction is already well handled by this model by using the zonal and 333 meridional component of the wind. So one big advantage of Random Forests 334 over linear regression is that it does not require to explore previously deeply 335 the data for extracting appropriate and relevant features as inputs to the 336

model. Figure 6 shows time series of 10m observed wind speed, NWP model 337 wind speed output over summer period of 2011 (panel a) and the probability 338 density function corresponding to the entire period, 2011 to 2015 (panel b). 339 Panels c and e show respectively time series of the reconstructed 10m wind 340 speed by LR_{SW}^{dir} (red line) and LR_{SW} (blue line), and by RF_A^{dir} (magenta 341 line) and RF_A (cyan line). Panels d and f show the corresponding probabil-342 ity density functions. All statistical models allows for a good bias correction. 343 All models underestimate the small quantiles of the distribution and give a 344 distribution very peaked around the mean observed wind speed. High per-345 centiles are however well reconstructed. This is encouraging because this 346 part of the distribution is important in terms of energy production. We can 347 nevertheless expect an overestimation of the wind energy production with 348 those models because of the underestimation of small percentiles. 349

100m wind speed reconstruction



Figure 7: Same as Figure 7, for 100m wind speed.

Figure 7 show results of the reconstruction of 100m wind speed with statistical models described in Table 5. LR_F allows a reduction of the RMSE of about 15% corresponding to 0.14 $m.s^{-1}$ and the best model LR_{SW}^{dir} reduces

15

the RMSE by 23% corresponding to 0.23 $m.s^{-1}$. The correlation is im-354 proved from 0.92 to 0.94. Adding the direction dependence to linear model 355 with only 100m wind speed (i.e. LR_F^{dir}) does not improve results regarding 356 LR_F . Indeed, the error on the 100m wind speed does not depend on the 357 direction. Using all explanatory variables (i.e. LR_A^{dir}) leads to a strong over-358 fitting. Surprisingly, the linear model using stepwise selection of explanatory 359 variables in each direction (i.e. LR_{SW}^{dir}) recovers an important information 360 as it performs significantly better than the other. Again, its adaptability 361 may be the cause of its good performance. The information on the direction 362 in Random Forests does not improve the results like for 10m wind speed 363 reconstruction. The more important variables for Random forests and step-364 wise choice are mainly the 100m wind speed, but also the wind shear in the 365 boundary layer. Lasso technique selects mainly 100m wind speed. 366



Figure 8: Same as Figure 6, for 100m wind speed.

Figure 8 shows time series of 100m observed wind speed, NWP model 367 wind speed output over summer period of 2011 (panel a) and the probabil-368 ity density function corresponding to the entire period from 2011 to 2015 369 (panel b). panel c and e show respectively time series of the reconstructed 370 100m wind speed by LR_{SW}^{dir} (red line) and LR_{SW} (blue line), and by RF_A^{dir} 37 (magenta line) and RF_A (cyan line). Panels d and f show the corresponding 372 probability density functions. Some peaked wind speeds are less overesti-373 mated after statistical downscaling. As for the 10m wind speed, statistical 374 models underestimate the small quantiles of the distribution and give a dis-375 tribution peaked around the mean observed wind speed. 376

377

To conclude, we built different statistical models to improve the representation of the 10m and 100m wind speed of the ECMWF analyses. It has been shown that the 100m wind speed in ECMWF analyses is already well represented as it displays no systematic bias and a good correlation. Never-

the the RMSE computed for the period 2011 to 2015 is still of $1.0 \ m.s^{-1}$. 382 Statistical models reduces the RMSE on the 10m wind speed between 40%383 and 65%, and between 15% and 23% for the 100m wind speed. They im-384 prove at the same time the correlation between the observed wind speed 385 and the reconstructed one. For linear models, the variables selection is of 386 great importance to avoid over-fitting, and an exploration step allows to im-387 proves results significantly. Random Forests give quite comparable results as 388 the best linear models, without needing variable selection and a preliminary 389 exploration of the data. 390

³⁹¹ 4 Forecasts of surface winds

In this section we use the previous statistical models based on the knowledge 392 of the observed wind speed and the outputs of ECMWF analyses to forecast 393 wind speed at five days horizon with a frequency of 3 hours. We have access 394 to 1 year of forecasts in 2015. We train these statistical models on ECMWF 395 analyses from 2011 to 2014, and apply the resulting model to the forecasts. 396 Figures 9 and 10 show respectively the RMSE averaged over the 365 forecasts 397 for the 10m and 100m wind speed. A strong diurnal cycle of the performances 398 of both ECMWF forecasts and downscaled statistical predictions of the 10m 399 wind speed is evidenced. This diurnal cycle seems to be observed also for 400 100m wind speed forecasts, but with a less important amplitude. As the 401 dataset is trained on the ECMWF analyses, we can affirm that diurnal cycle 402 is better represented in the ECMWF analyses than in ECMWF forecasts. 403 This could be indeed explained by the data assimilation system that may 404 help to correct errors coming from NWP model parameterizations. 405

An interesting result shown in Figure 9 is that performance of the LR_F 406 statistical model which is comparable to linear model LR_{SW} , showing that 407 the added value of other explanatory variables is valuable mainly for small 408 lead-times in this case. Adding the dependence with the direction (i.e. 409 LR_{SW}^{dir}) allows for a significant reduction of the RMSE. Random Forests 410 RF_A shows the best performance. In addition to the simplicity to fit this 411 model, its robustness seems to overcome linear regression models. For 100m 412 wind speed forecasts (Figure 10), Linear models LR_F , LR_{SW} , and LR_{SW}^{dir} 413 and Random Forest RF_A are comparable. 414

For both 10m and 100m wind speed forecasts, statistical downscaling models allow for significant improvements regarding ECMWF predicted wind speed, at any lead-time from 3 hours to 5 days. Training dataset on the analyses of ECMWF may not be optimal. Indeed, training a statistical model for each lead-time separately should deeply improve results. First, this could help to remove the displayed diurnal cycle, but may also let the increase in RMSE with the lead-time be less sharp.

422 5 Summary and concluding remarks

We have used statistical models to evaluate 10m and 100m wind speed at a given location from output of a NWP model. Comparison of the observed wind speed and ECMWF wind speed output at 10m and 100m within the 5 years of data show that ECMWF analyses well represent 100m wind speed.

17

17



Figure 9: RMSE, computed between the 10m observed wind speed, and the 10m forecast wind speed, averaged over the entire set of forecasts.



Figure 10: RMSE, computed between the 100m observed wind speed, and the 100m forecast wind speed, averaged over the entire set of forecasts.

The computed RMSE is of 1.0 $m.s^{-1}$ (the mean wind speed being of 5.8 $m.s^{-1}$) and no systematic bias is displayed. On the contrary, 10m wind speed output from ECMWF analyses displays a systematic overestimation the observed wind speed. The computed RMSE is of 1.4 $m.s^{-1}$ (the mean wind speed being of 2.4 $m.s^{-1}$).

By applying linear regression between a certain amount of selected vari-432 ables and observed wind speed, we reduce the RMSE for the 10m and 100m 433 reconstructed wind speed up to 65% and 23%, respectively. Those good 434 results have been achieved by fitting a linear model in 8 directions and by 435 automatic selection of valuable variables in those directions. Building such a 436 model thus requires a special treatment and a good knowledge of the specific 437 site so that it cannot be systematically applied to another site. Very inter-438 estingly, using Random Forests to reconstruct 10m and 100m wind speed 439 gives comparable results as the best linear models (about 57% and 20%, 440 respectively), while their performance is not sensitive to any preparation of 441 the data. Computing time is a bit longer than simple linear models, but it 442 is quite similar when a linear model is fitted in each direction. 443

In a second step, we applied the statistical models to forecast up to 5 days. Note that statistical models are trained on past analyses. Applying it on forecasts will work 'as well' only if the relationship between NWP

18

outputs and observed wind speed does not change with lead-time. This is not a-priori guaranteed as the analyses incorporate information from observation via data assimilation. Results are encouraging, because the RMSE between forecast wind speed and observed wind speed is significantly reduced compared to ECMWF forecasts whatever the lead-time, and for both 10m and 100m wind speeds. Interestingly, Random Forests perform the same or better than linear models when applied to forecast 10m or 100m wind speed.

The results obtained for the forecasts are very encouraging: even though 455 the training only involved analyses, the reduction in RMSE persisted for 456 lead-times up to 5 days. Promisingly, there are evident changes to be tried 457 which should lead to improvements of the performances. As a first, training 458 statistical downscaling models directly on ECMWF forecasts makes sense 459 as a transfer function adapted to each lead-time should take into account 460 the reduced performance of ECMWF forecasts around 15pm and thus cor-461 rect systematic errors in the representation of the diurnal cycle. Moreover, 462 training dataset for each lead-time separately should also help to reduce the 463 increase of RMSE with lead-time by adapting the explanatory variables to 464 forecasts performance. For instance, for short lead-time, statistical models 465 may find out that surface wind speed in ECMWF forecasts gather valuable 466 information so that this information would be used. It may nevertheless not 467 be the case at longer timescales, so that statistical models would prefer using 468 information from large-scale circulation (e.g pressure) which is well modeled 469 by ECMWF forecasts, even at lead-time up to 5 days. Secondly, the good 470 performance of Random Forests together with linear regression models de-471 notes that it is possible to reconstruct the wind speed with very different 472 relations. Model aggregation may thus be a path to retrieve more informa-473 tion than when using a single statistical model. 474

475

In this study, we choose to use only local informations coming from NWP outputs. Additive valuable informations may be retrieved from largerscale NWP outputs such as large-scale horizontal gradients of the pressure. However, the discussion on the added value of any other NWP outputs is site dependent, and is already part of research matters. For instance, it has been proved that large scale circulation patterns give valuable information at timescales up to months in some regions of France [Alonzo et al., 2017].

A wind farm is often equipped with many anemometers situated at 10m 483 and at the hub height, so that local intime observations are easily available as 484 well as wind power output. Forecasting wind speed using only NWP outputs 485 is a good way to improve forecasts, but local past observations may also be 486 used as explanatory variable. Indeed, at very short lead-time (3-hours), we 487 can assume that the error the NWP model make at t_{0h} (corresponding to 488 the analyses) may be correlated to the future error at time t_{3h} . We could 489 also imagine to create a direct link between NWP outputs and wind energy 490 production as in [Giorgi et al., 2011], using in addition the information on 491 the close past wind energy production at the condsidered wind farm. 492

As a preliminary illustration of the benefit of such an approach, we train Random Forests and Linear Regression with stepwise selection of variables to forecast 10m and 100m wind speed at time t_{3h} only, and add the error

19

19



Figure 11: RMSE computed over 10 k-fold forecasts of 10m (a) and 100m (b) wind speed at 3 hour lead-time, using the error on the 10m and 100m wind speed at time t_{0h} (denoted by $\Delta 10$ and $\Delta 100$, respectively) as an explanatory variable. The dashed line represent the averaged RMSE of Random Forest without using observations at t_{0h} , and boxes represents the RMSE over 10 k-fold forecasts.

on the wind speed at time t_{0h} as an explanatory variable of the future wind 496 speed at time t_{3h} . We perform a k-fold of 10 forecasts over the year 2015. 497 Results are represented in Figure 11. When forecasting 10m wind speed at 498 t_{3h} , using the error at time t_{0h} allows for a reduction of the RMSE of 5% 499 with Random Forests and of 10% with linear model compared to Random 500 Forest without the observation at time t_{0h} . When forecasting 100m wind 501 speed at t_{3h} , using the information on the 10m wind speed observed at t_{0h} 502 allows for an improvement of 2% to 6%. Adding the information on the 503 100m wind speed at time t_{0h} spectacularly improves results by 18% with 504 linear regression model. 505

506

In addition of the good results obtained when reconstructing 10m and 100m wind speed, we also showed encouraging results when forecasting wind speed up to 5 days. By using very different statistical models, we highlight advantages of Random Forests over multiple linear regression, e.g simplicity and robustness. Finally, very promising perspective for improving downscaling at short-term horizon is exposed ; it involves a pseudo-assimilation of a local past observed information into the statistical downscaling model.

514 Acknowledgements

This research was supported by the ANR project FOREWER (ANR-14-538 CE05- 0028). This work also contributes to TREND-X program on energy transition at Ecole Polytechnique. The authors thank Côme De Lassus 518 Saint Geniès and Medhi Kechiar who produced preliminary investigations 519 for this study.

520 **References**

- [RTE, 2016] (2016). Annual electricity report. Technical report, Reseau de
 Transport d'Electricite.
- [GWE, 2016] (2016). Global wind statistics. Technical report, Global Wind
 Energy Council.
- [Alonzo et al., 2017] Alonzo, B., Ringkjob, H., Jourdier, B., , Drobinski, P.,
 Plougonven, R., , and P.Tankov (2017). Modelling the variability of
 the wind energy resource on monthly and seasonal timescales. *Renewable Energy*, 113:1434–1446.
- ⁵²⁹ [Bauer et al., 2015] Bauer, P., Thorpe, A., and Brunet, G. (2015). "the quiet revolution of numerical weather prediction". *Nature*, 525:47–55.
- [Breiman, 1996] Breiman, L. (1996). Bagging predictors. In Machine Learn *ing*, pages 123–140.
- [Cassola and Burlando, 2012] Cassola, F. and Burlando, M. (2012). "wind
 speed and wind energy forecast through kalman filtering of numerical
 weather prediction model output". Applied Energy, 99:154–166.
- [Davy et al., 2010] Davy, R. J., Woods, M. J., Russell, C. J., and Coppin,
 P. A. (2010). Statistical downscaling of wind variability from meteorolog ical fields. *Boundary Layer Meteorol*, 175:161–175.
- ⁵³⁹ [Devis et al., 2013] Devis, A., van Lipzig, N., and Demuzere, M. (2013). A
 ⁵⁴⁰ new statistical approach to downscale wind speed distribution at a site in
 ⁵⁴¹ northern europe. Journal of geophysical research : Atmospheres, 118:2272–
 ⁵⁴² 2283.
- ⁵⁴³ [Folland et al., 2012] Folland, C. K., Scaife, A. A., Lindesay, J., and
 ⁵⁴⁴ Stephenson, D. B. (2012). "how potentially predictable is northern eu⁵⁴⁵ ropean winter climate a season ahead?". *Int. J. Climatol.*, 32:801–818.
- ⁵⁴⁶ [Friedman et al., 2001] Friedman, J., Hastie, T., and Tibshirani, R. (2001).
 ⁵⁴⁷ The elements of statistical learning, volume 1. Springer series in statistics
 ⁵⁴⁸ New York.
- ⁵⁴⁹ [Giorgi et al., 2011] Giorgi, M. G. D., Ficarella, A., and Tarantino, M.
 ⁵⁵⁰ (2011). "assessment of the benets of numerical weather predictions in
 ⁵⁵¹ wind power forecasting based on statistical methods". *Energy*, 36:3968–
 ⁵⁵² 3978.
- [Glahn and Lowry, 1972] Glahn, H. R. and Lowry, D. A. (1972). "the use of
 model output statistics (mos) in objective weather forecasting". Journal
 of Applied Meteorology, 11:1203–1211.
- ⁵⁵⁶ [Haeffelin et al., 2005] Haeffelin, M., Barthes, L., Bock, O., Boitel, C., Bony,
 ⁵⁵⁷ S., Bouniol, D., Chepfer, H., Chiriaco, M., Cuesta, J., Delanoe, J., Drobin⁵⁵⁸ ski, P., Dufresne, J.-L., Flamant, C., Grall, M., Hodzic, A., Hourdin, F.,
 ⁵⁵⁹ Lapouge, F., Lemaitre, Y., Mathieu, A., Morille, Y., Naud, C., Noel, V.,
 ⁵⁶⁰ OHirok, W., Pelon, J., Pietras, C., Protat, A., Romand, B., Scialom, G.,

21

and Vautard, R. (2005). "Sirta, a ground-based atmospheric observatory
 for cloud and aerosol research". Annales Geophysicae, 23:1–23.

- ⁵⁶³ [Howard and Clark, 2007] Howard, T. and Clark, P. (2007). "correction and
 ⁵⁶⁴ downscaling of nwp wind speed forecasts". *Meteorol. Appl.*, 14:105–116.
- ⁵⁶⁵ [Jung and Broadwater, 2014] Jung, J. and Broadwater, R. P. (2014). "cur⁵⁶⁶ rent status and future advances for wind speed and powerforecasting".
 ⁵⁶⁷ Renewable and Sustainable Energy Reviews, 31:762–777.
- ⁵⁶⁸ [Kalnay, 2003] Kalnay, E. (2003). Atmospheric Modeling, Data Assimilation
 ⁵⁶⁹ and Predictability. Cambridge University Press.
- ⁵⁷⁰ [Kubik et al., 2011] Kubik, M. L., Coker, P. J., and Hunt, C. (2011). Us⁵⁷¹ ing meteorological wind data to estimate turbine generation output: a
 ⁵⁷² sensitivity analysis. In *Proceedings of Renewable energy congress*, pages
 ⁵⁷³ 4074–4081.
- [Maraun et al., 2010] Maraun, D., Wetterhall, F., Ireson, M., Chandler, E.,
 Kendon, J., Widmann, M., Brienen, S., Rust, H. W., Sauter, T., Themel,
 M., Venema, V. K. C., Chun, K. P., Goodess, C. M., Jones, R. G., Onof,
 C., Vrac, M., and ThieleEich, I. (2010). "precipitation downscaling under
 climate change : Recent developments to bridge the gap between dynamical models and the end user.". *Review of Geophysics*, 48:RG3003.
- [Mohandes et al., 2011] Mohandes, M., Rehmanb, S., and Rahmand, S.
 (2011). "estimation of wind speed prole using adaptive neuro-fuzzy inference system (anfis)". Applied Energy, 88:4024–4032.
- [Salameh et al., 2009] Salameh, T., Drobinski, P., Vrac, M., and Naveau, P.
 (2009). "statistical downscaling of near-surface wind over complex terrain in southern france". *Meteorol Atmos Phys*, 103:253–265.
- [Soman et al., 2010] Soman, S. S., Zareipour, H., and O. Malik, P. M.
 (2010). "a review of wind power and wind speed forecasting methods with
 different time horizons". North American Power Symposium (NAPS),
 pages 1–8.
- [Tibshirani, 1994] Tibshirani, R. (1994). Regression shrinkage and selection
 via the lasso. Journal of the Royal Statistical Society, Series B, 58:267–
 288.
- [Vallis, 2006] Vallis, G. K. (2006). Atmospheric and Oceanic Fluid Dynam *ics.* Cambridge University Press, Cambridge, U.K.
- ⁵⁹⁵ [Wagenbrenner et al., 2016] Wagenbrenner, N. S., Forthofer, J. M., Lamb,
 ⁵⁹⁶ B. K., Shannon, K. S., and Butler, B. W. (2016). "downscaling surface
 ⁵⁹⁷ wind predictions from numerical weather prediction models in complex
 ⁵⁹⁸ terrain with windninja". Atmos. Chem. Phys., 16:5229–5241.
- [Wilby and Dawson, 2013] Wilby, R. L. and Dawson, C. W. (2013). "the
 statistical downscaling model: insights from one decade of application". *Int. J. Climatol.*, 33:1707–1719.

[Wilby et al., 1998] Wilby, R. L., Wigley, T. M. L., Conway, D., Jones, P. D.,
Hewitson, B., Main, J., and Wilks, D. S. (1998). "statistical downscaling
of general circulation model output: A comparison of methods". *Water Resources Research*, 34:2995–3008.

[Zamo et al., 2016] Zamo, M., Bel, L., Mestre, O., and Stein, J. (2016). "improved gridded wind speed forecasts by statistical postprocessing of numerical models with block regression.". Weather and Forecasting, 31:1929–1945.