

A non-parametric approach for wind speed distribution mapping

Freddy Houndekindo^{1*}, Taha B.M.J. Ouarda¹

¹Canada Research Chair in Statistical Hydro-Climatology, Institut national de la recherche
scientifique, Centre Eau Terre Environnement, 490 de la Couronne, Québec, QC, G1K 9A9,
Canada

*Corresponding author: Freddy Houndekindo
490, Couronne street, Québec, QC, G1K 9A9, Canada
Tel: +1 418-654-3842
E-mail: freddy.houndekindo@inrs.ca

21 Table of Contents

22	List of abbreviations	3
23	Highlights	4
24	Abstract	5
25	1. Introduction.....	6
26	2. Methodology	9
27	2.1. Quantile-based WS probability distribution mapping.....	10
28	2.1.1. MRMR approach for covariate selection.....	11
29	2.1.2. Regression models.....	12
30	2.1.3. Recovery of the WS distribution from WSQ.....	14
31	2.2. Weibull parameter mapping	16
32	2.3. Model validation	17
33	3. Study area and dataset.....	18
34	4. Results	20
35	4.1. Performance of regression models	20
36	4.2. Wind speed distribution mapping.....	23
37	5. Discussion	27
38	6. Conclusion	32
39	7. Acknowledgments	33
40	Appendix I Statistics of the estimated wind speed quantiles.....	34
41	Appendix II. Wind speed covariates	34
42	References.....	36

43

44

45

46

47

48

49

50

51 List of abbreviations

BS	Birnbaum-Saunders
CANGRD	Canadian gridded temperature and precipitation anomalies dataset
CDF	Cumulative probability function
D	Kolmogorov–Smirnov statistic
DEM	digital elevation model
ECDF	Empirical cumulative probability function
FS	Feature selection
GBT	Gradient boosting trees
GG	Generalized Gamma
KCDF	Kernel estimator of cumulative distribution function
LN	Log-Normal
LR	Linear regression
LSE	Least Square Estimation
MI	Mutual information
MRMR	Minimum redundancy maximum relevance
MRMR-MI	Minimum redundancy maximum relevance with Mutual information
MRMR-PC	Minimum redundancy maximum relevance with Pearson correlation coefficient
PC	Pearson correlation coefficient
PDF	Probability distribution function
PP plot	Percentage probability plot
QWSM	Quantile-based wind speed probability distribution mapping
R	Rayleigh
RD	Regional distribution
SSE	Sum of the square error
W	Weibull
WPM	Weibull parameters mapping
WS	Wind speed
WSQ	Wind speed quantile
XGB	Extreme Gradient Boosting

52

53

54

55

56

57

58

59 **Highlights**

- 60 • A non-parametric approach for wind speed mapping is developed.
- 61 • A comparative analysis of parametric and non-parametric approaches is carried out.
- 62 • The non-parametric method slightly outperforms the parametric approach and avoids the
- 63 hypothesis of a single distribution.
- 64 • The new method is recommended for regions having diverse wind regimes.

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84 **Abstract**

85 Statistical methods to estimate wind resources at unsampled locations in a region can serve as an initial
86 step to identify locations that warrant further investigation. There has been an ongoing effort to develop
87 approaches for mapping the parameters of the wind speed distribution with statistical methods. This
88 approach enables a comprehensive understanding of the wind resource variability across the entire
89 region by considering the full wind speed distribution rather than focusing solely on mean values. The
90 present study proposes a non-parametric approach to map the wind speed distribution. The method's
91 main advantage is that it avoids constraining the region to a single distribution family and is thus more
92 flexible than existing methods. In the proposed approach, a number of wind speed quantiles are first
93 mapped in the region using machine learning techniques. Afterwards, the wind speed distribution is
94 estimated by fitting an asymmetric kernel estimator to the estimated wind speed quantiles at unsampled
95 locations. The new approach was compared to the standard statistical method based on mapping the
96 regional wind speed distribution parameters. The results indicate that the non-parametric approach
97 leads in the best scenario to a 9% and 6% drop in the Kolmogorov-Smirnov statistic on average during
98 cross-validation and validation, respectively. The Birnbaum-Saunders and the Log-Normal kernels gave a
99 better fit to the estimated wind speed quantiles than the Weibull kernel. The proposed approach is
100 recommended in regions with high wind regime variability.

101 **Keywords:** Asymmetric kernel estimator, Non-parametric, Quantile, Wind speed distribution, wind
102 variability, ungauged location, regional estimation.

103

104

105

106 1. Introduction

107 Wind energy has the potential to become a crucial source of power worldwide [1]. In 2021, worldwide
108 wind energy installed capacity reached 837 GW, with an estimated offset of over 1.2 billion tons of CO₂
109 [2]. However, more effort is needed to raise the contribution of wind energy in the world energy mix to
110 achieve a more sustainable and low-carbon future [3].

111 One of the initial stages of building a wind farm involves finding a suitable location with sufficient wind
112 resources to generate electricity. This objective typically involves conducting an in-depth assessment of
113 the wind regime, which requires a long-term dataset of wind speed measurements. However, this data is
114 often only available at irregular points in space rather than at the location of interest for wind energy
115 production. It may not be feasible to install a monitoring station to gather sufficient data during the
116 preliminary site selection due to time and financial constraints. Using methods that can estimate wind
117 resources at unsampled locations is more suitable. Although these methods may not be as accurate as a
118 monitoring station, they can help identify potential sites that warrant further investigation.

119 Numerous WS estimation studies have been conducted at unsampled locations, as detailed in the review
120 by Houndekindo and Ouarda [4]. These studies typically estimate an aggregated WS value [5, 6], such as
121 the mean and occasionally the WS distribution, via mapping the parameters of a theoretical probability
122 distribution function. Both approaches have some downsides. First, using the mean WS for wind
123 resource assessment may underestimate the long-term resource depending on the frequency
124 distribution's shape [7]. Second, when estimating the WS distribution at unsampled locations, authors
125 typically select a unique family of distributions with different parameters for the entire region (the
126 regional distribution (RD)). For example, Veronesi, et al. [8] mapped WS distribution in the UK using
127 random forests and assumed that the Weibull distribution (W) was adequate across the study region.
128 Although the W is the most commonly used distribution for WS modelling, some studies have found that
129 other types of distributions may provide a better fit depending on the wind regime at a location. For

130 instance, the three-parameter W distribution (an additional location parameter) is better suited for calm
131 wind regimes [9]. Tsvetkova and Ouarda [10] reported that the heavy-tailed Halphen distribution family
132 provided a better fit than the two-parameter W distribution in all 125 WS stations considered in Eastern
133 Canada.

134 In another study, Jung [11] mapped WS distribution parameters in Southwest Germany. First, the author
135 evaluated the goodness of fit (GOF) of 67 theoretical distributions to select the RD. Then, a gradient-
136 boosting model was employed to map the parameters of the selected distribution. Similarly, Laib and
137 Kanevski [12] conducted a study in Switzerland for extreme WS. The authors used the quantiles plot to
138 evaluate the GOF of three theoretical distributions and select a RD. Then, with a machine learning
139 model, they mapped the parameters of the RD. This approach can be tedious, requiring the testing of
140 multiple distributions, and there is no guarantee that the selected distribution would be adequate at the
141 unsampled locations of interest. Previous studies evaluated the goodness of fit of different theoretical
142 distributions for WS modelling in a given region [13-18] and found that no single distribution family
143 provided the best fit at all locations in the region. Thus, using a single family of distributions may not be
144 appropriate for characterizing the WS distribution in an entire region.

145 This work proposes a new approach for WS distribution mapping that does not constrain the region to a
146 single distribution family (i.e.: a regional distribution). The proposed approach consists of estimating
147 several WS quantiles (WSQ) at a location of interest. Then, a distribution function can be fitted to the
148 estimated WS quantiles using the Least Square Estimation (LSE) method.

149 It can be tedious to test several distributions with the LSE method. Indeed, in most cases, the LSE
150 method does not have an analytical solution. Thus, optimization algorithms may be required with an
151 initial guess of the parameters, which can lead to suboptimal solutions. To address this issue, it is
152 proposed to fit a kernel estimator of cumulative distribution function (KCDF) to the estimated WSQ.

153 Kernel estimators are, in general, rather flexible and do not require prior knowledge of the family of
154 distributions of the data. The literature shows a growing interest in kernel estimators for WS distribution
155 modelling [19]. In most of these studies, symmetric kernels (ex: gaussian) were used to estimate the
156 probability distribution function. WS values are non-negative, while symmetric kernels have unbounded
157 support leading to probability leakage below zero [20]. This is a well-known problem called the boundary
158 effect, and several solutions have been proposed [21]. In this study, one of these solutions based on
159 asymmetric kernel estimators [22] is adopted and introduced for WS distribution modelling. According to
160 Hirukawa [22], asymmetric kernels are weight functions with support on the unit interval $[0, 1]$ or the
161 positive half-line. The effectiveness of the proposed approach was assessed by comparing it to another
162 method based on mapping the W parameters in the study region.

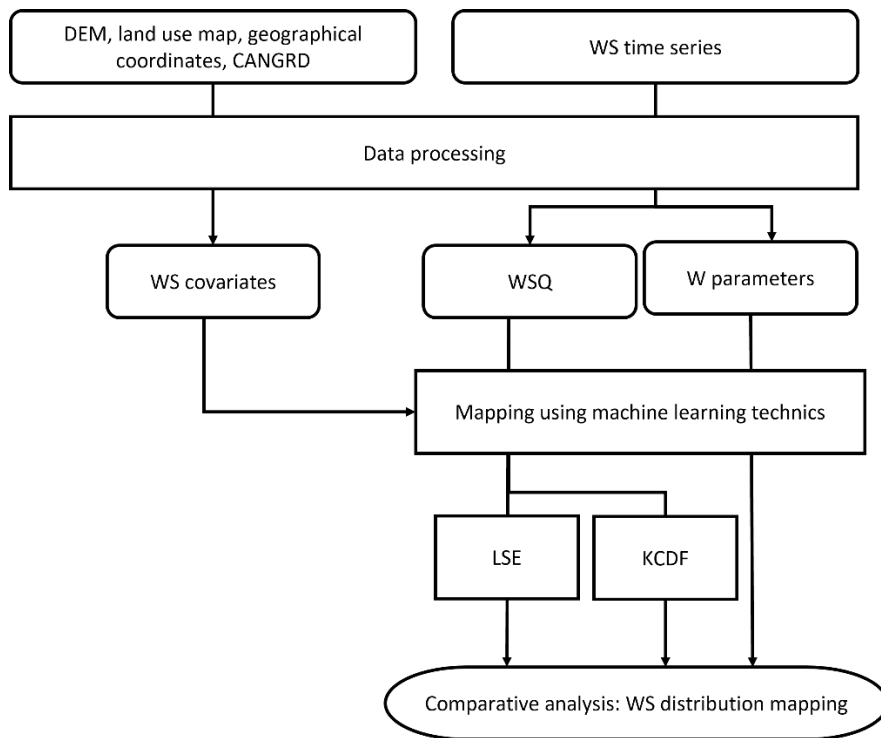
163 The paper's novelty can be summarized as follows: First, a methodology to map WS distribution is
164 proposed based on mapping WSQ. Quantiles are relatively easy to estimate from time series, while
165 selecting an adequate RD can be tedious, requiring the fitting and evaluation of multiple distributions.
166 Secondly, to the author's knowledge, this is the first study employing asymmetric kernels to model WS
167 distribution. By combining the mapping of WSQ and asymmetric kernels, a fully non-parametric
168 approach for WS distribution mapping is proposed in this study. The main advantage of the non-
169 parametric approach is that it does not require specifying a unique distribution family to the region of
170 interest. This allows to effectively combine all the available data in the region to build a more robust
171 model in case the region does not have a homogenous wind regime which can be described by a single
172 family of distribution functions.

173 The current paper is structured as follows. Section 2 illustrates the methodology of the proposed
174 approach with the evaluation procedure. The study area and the dataset are presented in section 3. The
175 results obtained are shown in section 4. In sections 5 and 6, the discussion of the findings and the
176 conclusion are given, respectively.

177 2. Methodology

178 This study proposes a new approach for mapping WS distribution using regional information without
179 constraining the region to a single distribution family. First, various WSQ are estimated at sampled
180 locations in the region. Then, machine learning and WS covariates are used to map the quantiles,
181 allowing the estimation of these WSQ at any unsampled location in the region. Finally, parametric, and
182 non-parametric approaches are implemented to recover the WS distribution at unsampled locations
183 from estimated quantiles. The proposed approach will be referred to as Quantile-based WS probability
184 distribution Mapping (QWSM) in the next sections. The QWSM approach will be compared to another
185 approach based on directly mapping the W parameters [8]. This method will be referred to as the W
186 parameters mapping (WPM) in the next sections. A flowchart of the methodology is available in Figure 1.

187



188

189 Figure 1: Methodology of the comparative analysis of WS probability distribution mapping approaches

190 2.1. Quantile-based WS probability distribution mapping

191 At the sampled locations in the region, WSQ at some fixed percentile points can be estimated from the
192 sorted values of the hourly time series with the following general formula [23]:

193
$$W(P) = (1 - \gamma)X_{(j)} + \gamma X_{(j+1)} \tag{1}$$

194 Where P is the percentile point of interest, $X_{(j)}$ and $X_{(j+1)}$ are j -th order statistics. γ is a weight ($0 \leq$
195 $\gamma \leq 1$) that is function of $j = \text{floor}(Pn + m)$, $m = \alpha + P(1 - \alpha - \beta)$ and $g = nP + m - j$. In case it
196 is desired to obtain $W(P)$ as a continuous function of P , then $\gamma = g$ and selecting γ reduces to selecting
197 α, β . Typical values of α, β are available in [23]. In this study, α, β were both set to $1/3$ given quantiles
198 that are approximately median-unbiased regardless of the WS true probability distribution [24]. Using
199 equation 1, WSQ associated with the following 13 percentile points were estimated at the sampled
200 locations: 5.0% (P1), 12.5% (P2), 20.0% (P3), 27.5% (P4), 35.0% (P5), 42.5% (P6), 50.0% (P7), 57.5.0% (P8),
201 65.0% (P9), 72.5% (P10), 80.0% (P11), 87.5% (P12), and 95.0% (P13). Table 6 in Appendix I gives an
202 overview of the distribution of the estimated WSQ.

203 These percentile points were chosen to cover the WS cumulative distribution functions (CDF) evenly,
204 ensuring a representative estimation of the WSQ at various points along the distribution. In previous
205 studies employing a similar modelling approach, varying numbers of percentile points have been
206 modelled to estimate the probability distribution of a target variable. For instance, to forecast power
207 load probability distribution, [25] modelled 20 percentiles evenly spaced between 1% and 96%. In
208 another study, to map wind speed shear distribution, [26] estimated 11 percentiles evenly spaced
209 between 1% and 99%. Additionally, to regionalize river temperature at ungauged locations, [27]
210 estimated 17 percentiles non-evenly spaced between 0.05% and 99.95%. This diversity in the number of
211 percentile point selections highlights a lack of consensus in the literature regarding the optimal number

212 to ensure a comprehensive target distribution coverage. Nevertheless, it is worth noting that the number
213 of percentiles selected in the current study falls within the range of those used in previous research.
214 A regression function was constructed between the observed WSQ and WS covariates. Two regression
215 models were compared, the multilinear regression (LR) and the Gradient boosting trees [GBT: 28] model.
216 Feature selection (FS) was performed using the minimum redundancy maximum relevance (MRMR)
217 method [29] to reduce the complexity of the models and improve their performance. A comparative
218 study of FS methods was carried out by Houndekindo and Ouarda [30]. They found that MRMR was
219 among the most effective FS methods for WSQ estimation. Houndekindo and Ouarda [30] used MRMR
220 with simple linear regression. However, the approach can be adapted to non-linear models such as tree-
221 based gradient boosting. The FS method (MRMR) and the GBT model are presented in more detail in the
222 following subsections.

2.1.1. MRMR approach for covariate selection

224 MRMR is a filter-based FS approach with the benefit of considering both the covariates' relevancy and
225 redundancy during selection. Filter-based FS methods are computationally efficient algorithms and are
226 agnostic to the regression model [31]. The MRMR algorithm uses an iterative approach to select the
227 covariate (X_i) at each step with the best trade-off between its relevancy to the response variable (Y) and
228 its redundancy relative to selected features from previous iterations. At the first step of the algorithm,
229 the most relevant covariate is selected based on a measure of relevancy ($Rel(X_i, Y)$).

230 Let $Red(X_i, X_j)$ be a measure of the dependency between the covariates X_i and X_j and let S be the set
231 of covariates selected during previous iterations. After the first step of the algorithm, S contains only the
232 most relevant covariate ($\max_{X_i} [Rel(X_i, Y)]$) and the objective criterion at each subsequent iteration of
233 the MRMR algorithm can be formulated in two ways:

234 $\max_{X_i \notin S} [Rel(X_i, Y) / Red(X_i, X_j)]$ (2)

235 Or

236 $\max_{X_i \notin S} [Rel(X_i, Y) - Red(X_i, X_j)]$ (3)

237 Several measures of relevancy and redundancy can be applied. In this study the following formulations
 238 of the MRMR objective criterion were compared:

239 $MRMR - PC: \max_{X_i \notin S} \left[F(X_i, Y) / \left(\frac{1}{S} \sum_{X_j \in S} \rho(X_i, X_j) \right) \right]$ (4)

240 and

241 $MRMR - MI: \max_{X_i \notin S} \left[I(X_i, Y) / \left(\frac{1}{S} \sum_{X_j \in S} I(X_i, X_j) \right) \right]$ (5)

242 Where $F(X_i, Y)$ is the F-statistic used to measure the relevancy, $\rho(X_i, X_j)$ is the Pearson correlation
 243 coefficient (PC) used to measure redundancy, $I(X_i, Y)$ is the mutual information (MI) used to measure
 244 relevancy and $I(X_i, X_j)$ is the MI used to measure redundancy. The MI between two random variables X
 245 and Y can be defined as follows:

246 $I(X, Y) = \iint p(X, Y) \log(p(X, Y) / p(X)p(Y)) \, dx dy$

247 (6)

248 The Python package scikit-learn [32] was used to calculate the MI between the variables.

249 2.1.2. Regression models

250 The LR model was implemented and used as a benchmark for the GBT model. Tree-based regression
 251 models such as GBT perform better than deep learning models on tabular data and often outperform
 252 other regression models [33]. The GBT algorithm works by fitting sequentially decision trees to the
 253 residuals from previous iterations. Contrary to the LR model, the GBT model can learn nonlinear

254 relationships between the covariates and the response variable and is robust against non-informative
 255 covariates [34]. The GBT model is a popular regression model that has been successfully applied in
 256 studies for short-term wind power prediction [35], wind resource mapping [26], the selection of solar
 257 power plant location [36] and short-term prediction of solar irradiance [37].

258 The eXtreme Gradient Boosting package [XGB: 38] is a popular machine-learning library that implements
 259 the GBT algorithm efficiently. Several regularization strategies are available in XGB to improve the model
 260 performance and reduce computational time. To find adequate values for the parameters of XGB, a
 261 random search with 1000 iterations was implemented. Grid search and random search are popular
 262 algorithms used for hyperparameter tuning [39]. Grid search is a brute force algorithm that
 263 systematically tries all possible combinations of hyperparameter values within specified ranges. The
 264 algorithm can find the optimal hyperparameter values within the defined search space at the cost of
 265 increased computational resources and time. On the other hand, random search is a more efficient
 266 algorithm that does not guarantee the optimal solution but can find good hyperparameters [40]. Table 1
 267 presents the hyperparameters of the XGB model that were tuned in the study.

268 Table 1: Hyperparameters of the XGB model

Hyperparameters used during training	Search space (Min, Max, Step)
Learning rate (Boosting learning rate)	(0.01, 0.1, 0.01)
Minimum loss reduction (gamma)	(0.0, 1.0, 0.1)
Maximum depth of the trees (max_depth)	(3, 10, 1)
Ratio of predictor to use during training (colsample_bytree)	(0.1, 0.7, 0.1)
Subsample ratio of the training data (subsample)	(0.1, 0.5, 0.1)

269

270 2.1.3. Recovery of the WS distribution from WSQ

271 With estimated WSQ available at any non-sampled location, it is possible to fit different theoretical
 272 distribution functions using the LSE method. The LSE method is widely used for fitting WS probability
 273 distributions [41]. In their study, Jung and Schindler [26] applied the LSE method to recover the
 274 probability distribution of wind shear exponent from estimated quantiles of the same variable. LSE
 275 involves minimizing the sum of the square error (SSE) between the empirical cumulative probability
 276 (ECDF) and the theoretical CDF to determine the best-fitting parameters of the theoretical distribution
 277 function. Let \widehat{W}_i be the predicted WSQ and $\widehat{F}(W_i)$ their associated CDF, the SSE can be written as
 278 follows:

$$279 \quad SSE = \sum_{i=1}^{13} [\widehat{F}(W_i) - F(\widehat{W}_i; \hat{\theta})]^2 \quad (7)$$

280 Where: $F(\widehat{W}_i; \hat{\theta})$ corresponds to the cumulative probability function of \widehat{W}_i with estimated parameter $\hat{\theta}$.
 281 The W, Log-Normal (LN), Rayleigh (R) and Generalized Gamma (GG) distribution were fitted to the
 282 estimated WSQ.

283 Additionally, it is proposed to recover the WS distribution at unsampled locations using asymmetric
 284 KCDF. The asymmetric kernels method represents one of the solutions to the boundary effects that
 285 appear when using symmetric kernels with bounded random variables (ex.: WS values are bounded on
 286 $[0, \infty]$). By combining WSQ mapping and asymmetric kernel fitting, this study proposes a fully non-
 287 parametric method for wind speed distribution mapping. Traditional parametric methods might
 288 introduce bias if the selected RD does not align with the data. The non-parametric approach can adapt to
 289 various WS distribution patterns without being restricted by specific parametric assumptions. This
 290 flexibility is necessary for a region with complex and diverse wind behaviors. In addition, combining the

291 WSQ mapping and asymmetric kernel fitting avoids the tedious process of testing and evaluating
 292 different probability distribution functions to model WS.

293 The general expression for the asymmetric KCDF is given by [21]:

$$294 \hat{F}(w) = 1/n \sum_{i=1}^n \bar{K}_{w,b}(W_i),$$

295 (8)

296 Where:

297 $b > 0$ is the bandwidth and $\bar{K}(\cdot)$ is the CDF of an asymmetric kernel function. In this work, the
 298 Birnbaum-Saunders (BS), the Log-Normal (LN) and W asymmetric kernel functions were tested [21, 42]:

$$299 \hat{F}^{BS}(w) = 1/n \sum_{i=1}^n \bar{K}_{BS}(W_i; w, \sqrt{b}), \quad (9)$$

$$300 \hat{F}^{LN}(w) = 1/n \sum_{i=1}^n \bar{K}_{LN}(W_i; \log w, \sqrt{b}), \quad (10)$$

$$301 \hat{F}^{WB}(w) = 1/n \sum_{i=1}^n \bar{K}_{WB}(W_i; w/\Gamma(1+b), 1/b), \quad (11)$$

302 Where:

$$303 \bar{K}_{BS}(x; \beta, \alpha) = 1 - \Phi((\sqrt{x/\beta} - \sqrt{\beta/x})/\alpha), \quad \beta, \alpha > 0,$$

304 (12)

$$305 \bar{K}_{LN}(x; \mu, \sigma) = 1 - \Phi((\log x - \mu)/\sigma), \quad \mu, \sigma > 0,$$

306 (13)

$$307 \bar{K}_{WB}(x; \alpha, \beta) = \exp(-(x/\beta)^\alpha), \quad \alpha, \beta > 0, \quad (14)$$

308 $\Phi(\cdot)$ is the CDF of the standard normal distribution and $\Gamma(\cdot)$ is the gamma function.

309 The optimal bandwidths can be selected by minimizing the Mean Integrated Square Error (MISE)

$$310 MISE = \int_0^\infty MSE(\hat{F}(w)) dw \quad (15)$$

311 Where:

$$312 \quad MSE(\hat{F}(w)) = E \left[\left(\hat{F}(w) - F(w) \right)^2 \right] \quad (16)$$

313 Mombeni, et al. [21] derived the asymptotical optimal bandwidth of \bar{K}_{BS} and \bar{K}_{WB} with respect to the

314 MISE:

$$315 \quad b_{opt}^{BS} \approx \left\{ \int_0^\infty x f(x) dx \right\}^{2/3} \left\{ \pi^2 \int_0^\infty (x f(x) + x^2 f'(x))^2 dx \right\}^{-2/3} n^{-2/3},$$

316 (17)

$$317 \quad b_{opt}^{WB} \approx \left\{ 36 \ln 2 \int_0^\infty x f(x) dx \right\}^{1/3} \left\{ \pi^4 \int_0^\infty (x^2 f'(x))^2 dx \right\}^{-1/3} n^{-1/3}, \quad (18)$$

318 Lafaye de Micheaux and Ouimet [42] proposed the following asymptotical optimal bandwidth with

319 respect to the MISE for \bar{K}_{LN} :

$$320 \quad b_{opt}^{LN} \approx \left\{ \frac{1}{\sqrt{\pi}} \int_0^\infty x f(x) dx \right\}^{2/3} \left\{ 4 \int_0^\infty \frac{x^2}{4} (f(x) + x f'(x))^2 dx \right\}^{-2/3} n^{-2/3}, \quad (19)$$

321 The optimal bandwidth with respect to the MISE was selected under the assumption that the W with

322 parameters estimated using the predicted WSQ and the LSE method was the target distribution. The

323 reason for employing the W distribution in the paper is two-fold: First, it is the parametric probability

324 distribution function most commonly used to model WS; Secondly, it is convenient because its CDF can

325 be linearized with respect to its parameters and the WSQ. As a result, finding the best-fitting parameters

326 with the LSE method is equivalent to solving a linear equation and does not require an optimization

327 algorithm.

328 2.2. Weibull parameter mapping

329 In previous studies, to estimate the WS probability distribution at unsampled locations, machine learning

330 models were used to map the parameters of a RD. The approach selects a single distribution family for

331 the entire region. Then, the distribution function parameters are fitted at the sampled locations, and a
332 regression model is built between the parameters and WS covariates. Jung [11] selected the Wakeby
333 distribution as the RD in southwest Germany based on two goodness of fit measures: Kolmogorov-
334 Smirnov statistic and the coefficient of determination. For a review of criteria used for the identification
335 of adequate WS distributions the reader is referred to Ouarda, et al. [43]. Veronesi, et al. [8] selected the
336 W as the RD in the UK due to its widespread use in modelling WS, and convenience as it requires only
337 two parameters to characterize the WS probability distribution. The W was also adopted as the RD in this
338 study to evaluate the QWSM approach. The W parameters were estimated with the LSE method and the
339 best-fitting parameters were mapped in the region using the WS covariates described in section 3 and
340 the LR and XGB regression models described in section 2.1.2. The MRMR algorithm was also applied to
341 identify the best set of covariates to include in the regression models.

342 2.3. Model validation

343 To evaluate the QWSM and the WPM, holdout and 5-fold cross-validation were implemented with the
344 available samples. During the holdout procedure, parts of the samples were withheld (the validation set)
345 before model training and parameter tuning and used to evaluate the final model generalization
346 performance. During 5-fold cross-validation, the training samples were divided into five approximately
347 equal subsets. Then, the holdout method was implemented five times by considering each subset as the
348 validation set and training the model on the remaining subsets.

349 The following metrics were calculated based on the observed (y_i) and estimated (\hat{y}_i) values:

$$350 R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (20)$$

$$351 RMSE = \sqrt{1/n \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (21)$$

$$352 MAE = 1/n \sum_{i=1}^n |y_i - \hat{y}_i| \quad (22)$$

353

354 The evaluation of the GOF of the estimated WS probability distribution was based on the percentage
355 probability plot [PP plot: 44]. The PP plot compares the ECDF to the estimated CDF. During cross-
356 validation and validation, the R^2 , the RMSE and the MAE defined in equations 20, 21 and 22, respectively,
357 were used to evaluate the degree of association between the ECDF and the CDF. Horst [45] noted that
358 the PP plot has strong discriminatory power in high-density regions of the distribution (i.e.: the middle of
359 a distribution), where the CDF changes more rapidly with the WS values compared to low-density
360 regions (i.e.: the tails). Regions of the probability distribution with high density are the most crucial for
361 wind energy production. Also, in their reviews on WS distribution selection, Jung and Schindler [9]
362 observed that the most widely used GOF metrics were based on the PP plot.

363 The Kolmogorov–Smirnov statistic (D) is an alternative measure that was used to compare the ECDF and
364 the CDF:

$$D = \max |F_n(W_i) - \hat{F}(W_i)|$$

366 (23)

367 Where $F_n(W_i)$ is the ECDF and $\hat{F}(W_i)$ is the estimated CDF.

368 The ECDF was calculated with the Weibull plotting position [46] giving unbiased non-exceedance
369 probabilities regardless of the underlying distribution of the data [47]:

$$F_n(W_i) = i/(n + 1)$$

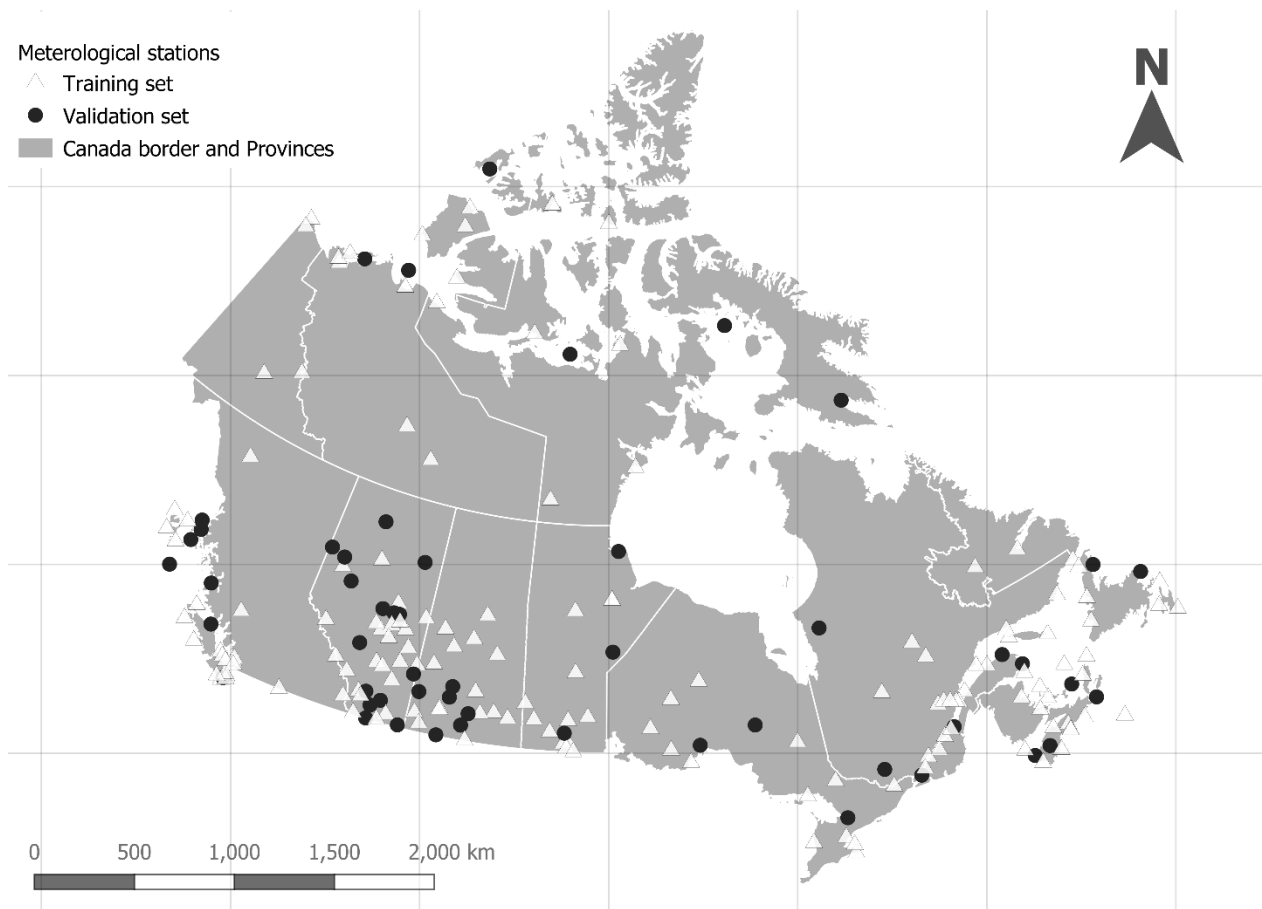
370 (24)

371 Where $i = 1, \dots, n$ is the rank of the WS values after sorting them in ascending order.

372 3. Study area and dataset

373 The study was conducted on data from the whole Canada representing a total area of 9,984,670 square
374 kilometers. Hourly WS data from 207 meteorological stations located throughout the country were used
375 for the research. From Environment and Climate Change Canada (ECCC) historical climate database,

376 stations with at least 20 years of recent WS record were selected. Additional filtering was performed to
377 eliminate all stations with more than ten years of record having two months of missing data. Figure 2
378 illustrates the geographical location of the 207 stations that were selected after filtering. From the
379 available stations, 155 (white triangles in figure 2) were used for FS, model training and cross-validation
380 and the remaining stations (black dots in figure 2) were used to validate the final model as explained in
381 section 2.3.



382

383 Figure 2: Spatial distribution of the training and validation stations used in this study.

384 The following four types of covariates were used with the regression models to either estimate the WSQ
385 or the W parameters: topographic, climatic, geographic, and surface roughness length. The
386 topographical covariates were created using the WhiteboxTools [48] and a 30m resolution global DEM

387 [49]. Seasonal and annual trends of mean temperature data were acquired from the Canadian gridded
388 temperature and precipitation anomalies (CANGRD) dataset (available at [https://climate-](https://climate-change.canada.ca/climate-data/#/historical-gridded-data)
389 [change.canada.ca/climate-data/#/historical-gridded-data](https://climate-change.canada.ca/climate-data/#/historical-gridded-data)). Surface roughness length was extracted from
390 a 2015 Canada land use map [50] resampled at different spatial resolutions using majority resampling
391 (i.e.: most popular value in a defined radius). A surface roughness length was associated with each land
392 use type based on a lookup table proposed by Wiernga [51]. Table 7 in Appendix II provides more details
393 about the covariates.

394 4. Results

395 4.1. Performance of regression models

396 The LR and the XGB models were fitted with covariates selected using MRMR-PC and MRMR-MI. The
397 results of comparing the different combinations of regression models and FS methods are presented in
398 Tables 3 and 4 for QWSM and the WPM, respectively. Figure 3 details the average R^2 for estimating the
399 13 WSQ and the two W parameters (shape and scale). The comparisons using cross-validation and
400 validation lead to very similar results, indicating, in general, that XGB with MRMR-PC outperforms the
401 other combinations of regression models and FS methods. Indeed, XGB gave better results than LR in
402 most cases, and MRMR-PC was more effective than MRMR-MI for FS in the study. In the few cases where
403 LR outperformed XGB, the performance difference was marginal and inconsistent during cross-validation
404 and validation (see, for instance, P8 in Figures 3a and 3b). Tables 3 and 4 indicate that the improved
405 performance of XGB with MRMR-PC is consistent across all metrics. Hereon, only the results obtained
406 with estimations from the top-performing FS and regression model (MRMR-PC + XGB) will be presented.
407 Figure 4 displays the spatial distribution of the RMSE (WSQ) scaled by the actual WS median for the
408 validation set. This representation allows for comprehensive visualization of the accuracy and variability
409 of the model's predictions across different locations. Scaling the RMSE with the actual median provides a

410 relative measure of error that can be compared and interpreted meaningfully. The spatial distribution of
 411 the scaled RMSE revealed that the model exhibited acceptable performances in estimating the WSQ in
 412 regions with sparse training samples highlighting its generalization capability.

413 Table 3: Average performance metrics for the estimation of WSQ

Validation Methods	Regression model	MRMR	MAE	R ²	RMSE
			km/h		km/h
Cross-validation	LR	MI	3.59	0.23	4.90
Cross-validation	LR	PC	3.40	0.26	6.11
Cross-validation	XGB	MI	3.24	0.42	4.30
Cross-validation	XGB	PC	3.08	0.47	4.07
Validation	LR	MI	3.64	0.36	4.48
Validation	LR	PC	3.24	0.46	4.19
Validation	XGB	MI	3.30	0.46	4.22
Validation	XGB	PC	3.00	0.57	3.74

414

415 Table 4: Average performance metrics for the estimation of the W parameters

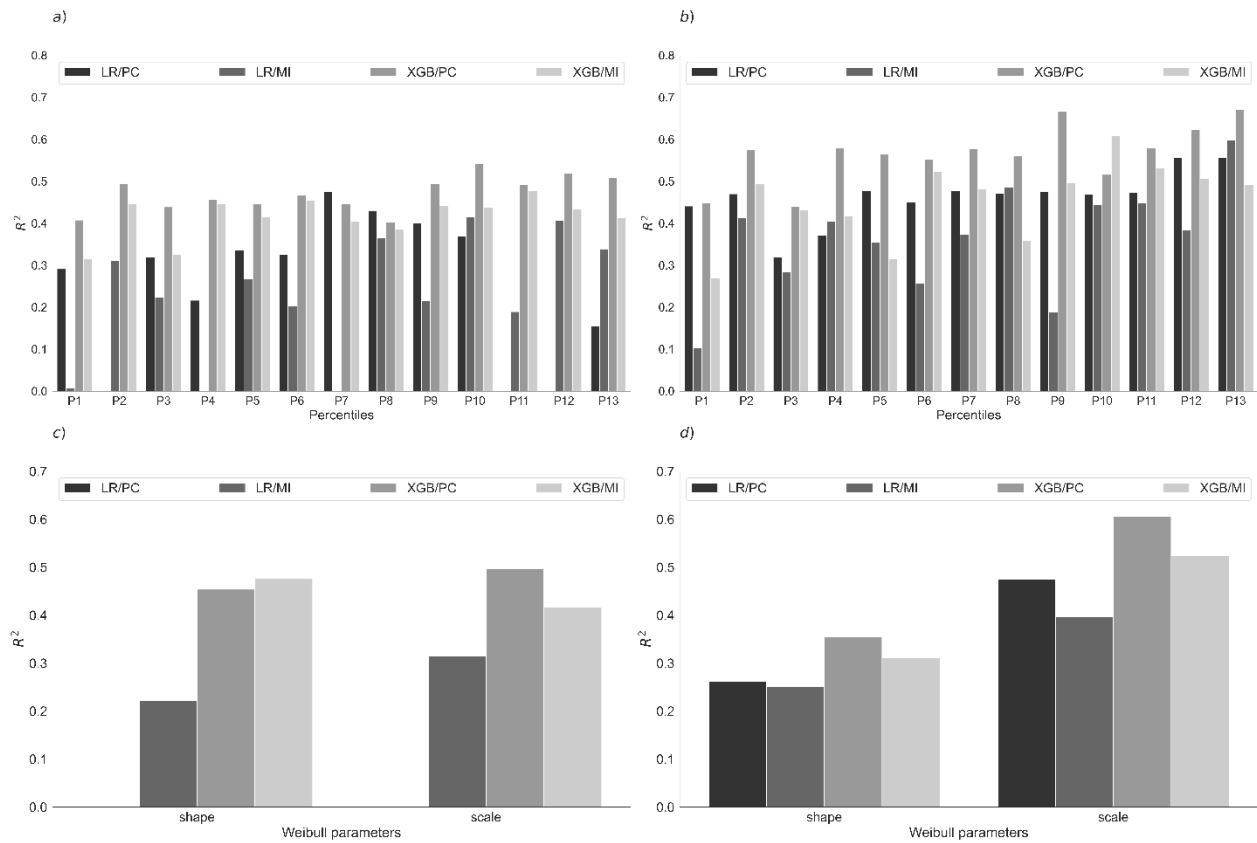
Validation Methods	Regression model	MRMR	MAE	R ²	RMSE
Cross-validation	LR	MI	1.88	0.27	2.47
Cross-validation	LR	PC	2.02	-	4.79
Cross-validation	XGB	MI	1.83	0.45	2.27
Cross-validation	XGB	PC	1.61	0.48	2.12
Validation	LR	MI	2.07	0.32	2.42
Validation	LR	PC	1.76	0.37	2.27
Validation	XGB	MI	1.75	0.42	2.16
Validation	XGB	PC	1.58	0.48	1.97

416

417

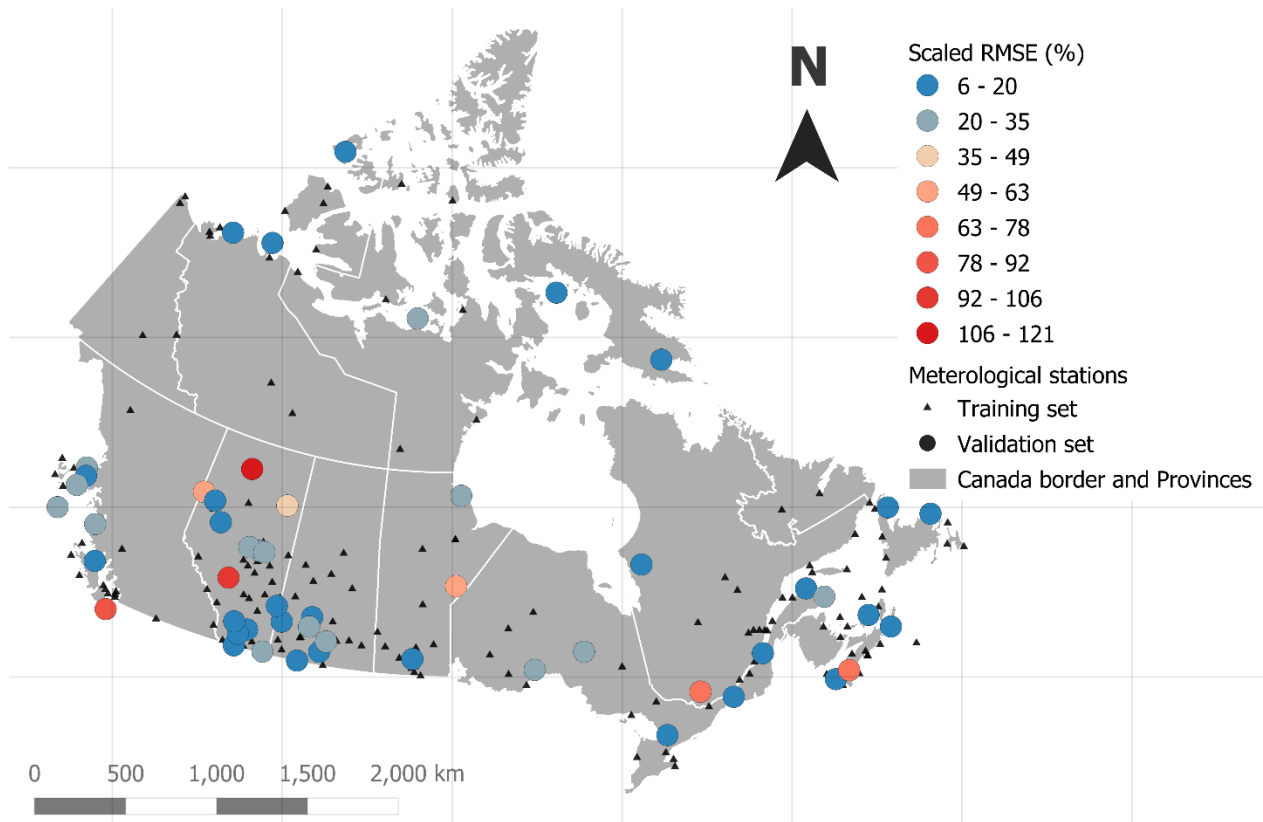
418

419



420

421 Figure 3: Performance of LR and XGB for the estimation of the WSQ (a and b) and the W parameters (c
 422 and d) during cross-validation (a and c) and validation (b and d). Note: Negative values of R^2 were set to
 423 zero



424
 425 Figure 4: Spatial distribution of the scaled RMSE (WSQ) of the validation set

426
 427 **4.2. Wind speed distribution mapping**

428 This section presents the results of the comparative analysis between the QWSM and WPM. Table 5
 429 shows the mean values of the GOF metrics. In general, it is observed that the QWSM gave a better fit
 430 than WPM for the considered metrics. Also, QWSM/W gave better fit than WPM. According to the R^2 ,
 431 RMSE and MAE criteria, QWSM/W and QWSM/GG were the best-performing methods, and their
 432 performances are very similar to QWSM/KCDF/BS and QWSM/KCDF/LN. However, during cross-
 433 validation and validation, the Kolmogorov-Smirnov statistic (D) seemed to favor QWSM/KCDF/LN and
 434 QWSM/KCDF/BS. The distribution of the GOF measures was represented using boxplots in Figure 4. The
 435 most noticeable difference in the distribution of the GOF measures was observed with D when
 436 comparing the different approaches. The methods based on QWSM/KCDF/LN and QWSM/KCDF/BS
 437 resulted in smaller D values and less variability in the same GOF measure compared to other approaches.

438 Furthermore, the different methods were evaluated by comparing the observed and estimated WSQ
 439 across ten equidistant percentiles ranging from 0.1 to 0.9. The outcome of this analysis (Figure 6)
 440 indicated that the QWSM methods often outperformed the WPM for the considered WSQ. Methods
 441 based on QWSM with the asymmetric kernels tend to give comparable performances to the parametric
 442 methods in the middle of the distribution (ex.: 0.4, 0.5, 0.6 percentiles). While in the tails (ex.:
 443 percentiles 0.1, 0.9) the parametric methods showcased a better performance than the non-parametric
 444 methods.

445

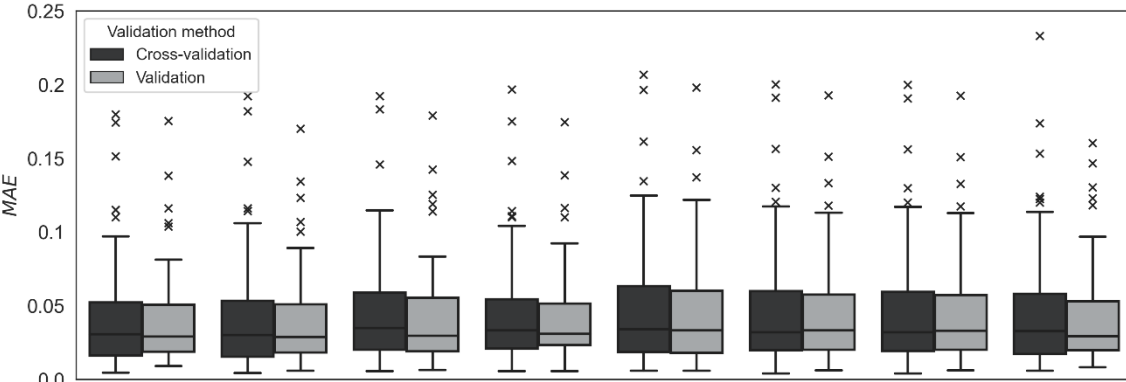
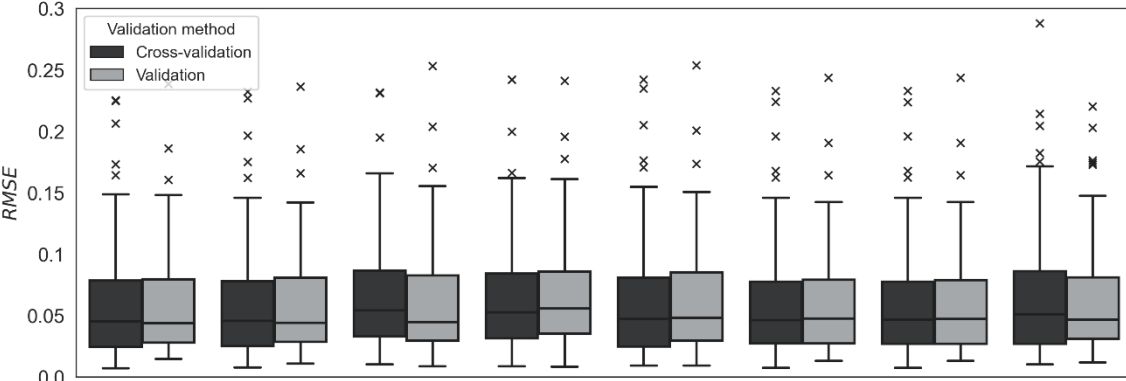
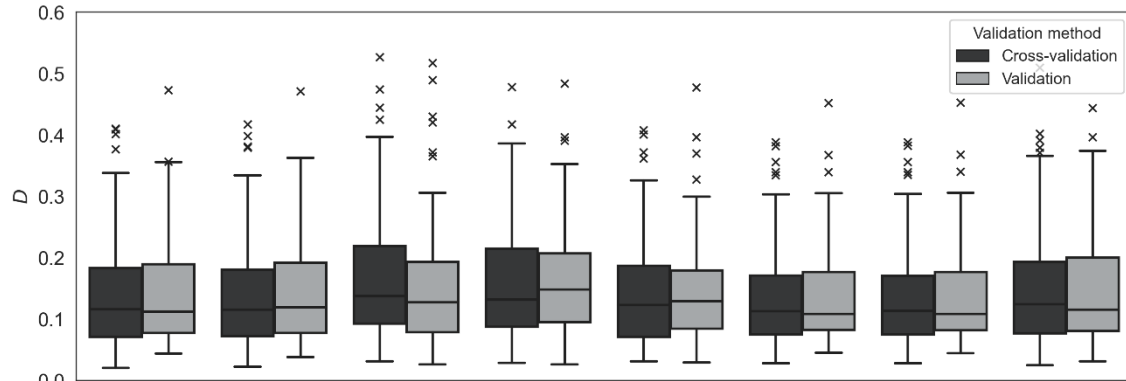
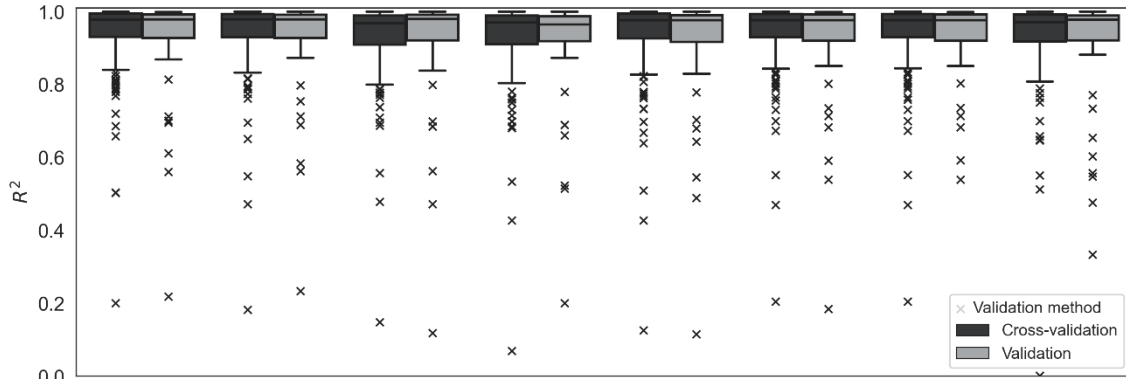
446 Table 5: Mean value of the GOF measures

Distribution	Validation Methods	D	MAE	R ²	RMSE
QWSM/GG	Cross-validation	0.137	0.039	0.938	0.058
QWSM/GG	Validation	0.147	0.041*	0.922*	0.062*
QWSM/KCDF/BS	Cross-validation	0.131	0.043	0.938	0.059
QWSM/KCDF/BS	Validation	0.143*	0.045	0.920	0.063
QWSM/KCDF/LN	Cross-validation	0.131	0.044	0.937	0.059
QWSM/KCDF/LN	Validation	0.143*	0.045	0.920	0.063
QWSM/KCDF/W	Cross-validation	0.137	0.046	0.932	0.061
QWSM/KCDF/W	Validation	0.150	0.046	0.911	0.064
QWSM/LN	Cross-validation	0.165	0.042	0.93	0.064
QWSM/LN	Validation	0.165	0.043	0.913	0.065
QWSM/R	Cross-validation	0.157	0.042	0.926	0.065
QWSM/R	Validation	0.168	0.044	0.908	0.069
QWSM/W	Cross-validation	0.136	0.039	0.939	0.058
QWSM/W	Validation	0.147	0.041*	0.921	0.062*
WPM	Cross-validation	0.144	0.042	0.93	0.062
WPM	Validation	0.152	0.043	0.910	0.065

Note: The best-performing methods are indicated in bold for the cross-validation and marked with * for the validation.

447

448 Table 5

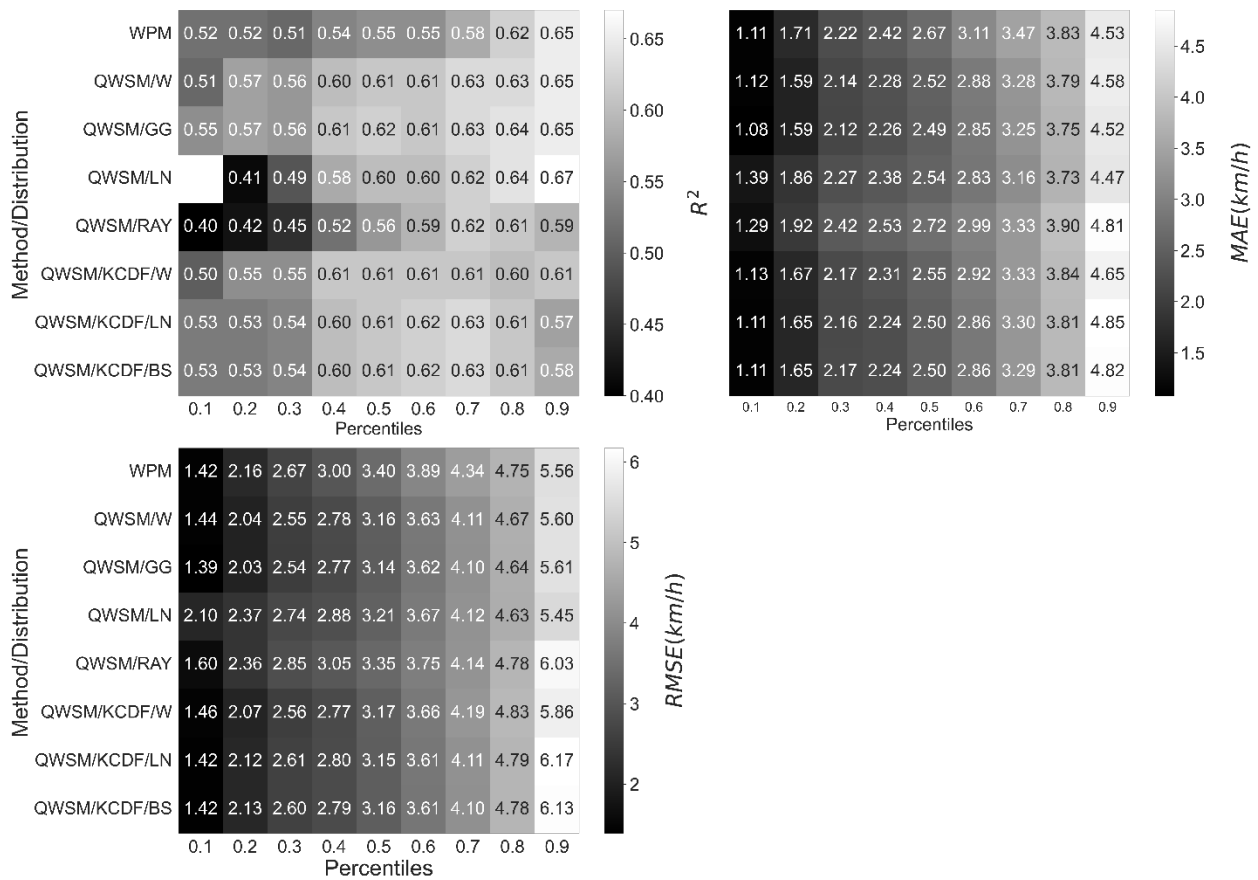


QWSM/W QWSM/GG QWSM/LN QWSM/R QWSM/KCDF/W QWSM/KCDF/LN QWSM/KCDF/BS WPM

Method/Distribution

450 Figure 5: GOF of estimated WS probability distribution. Note: Negative values of R^2 were set to zero

451



452

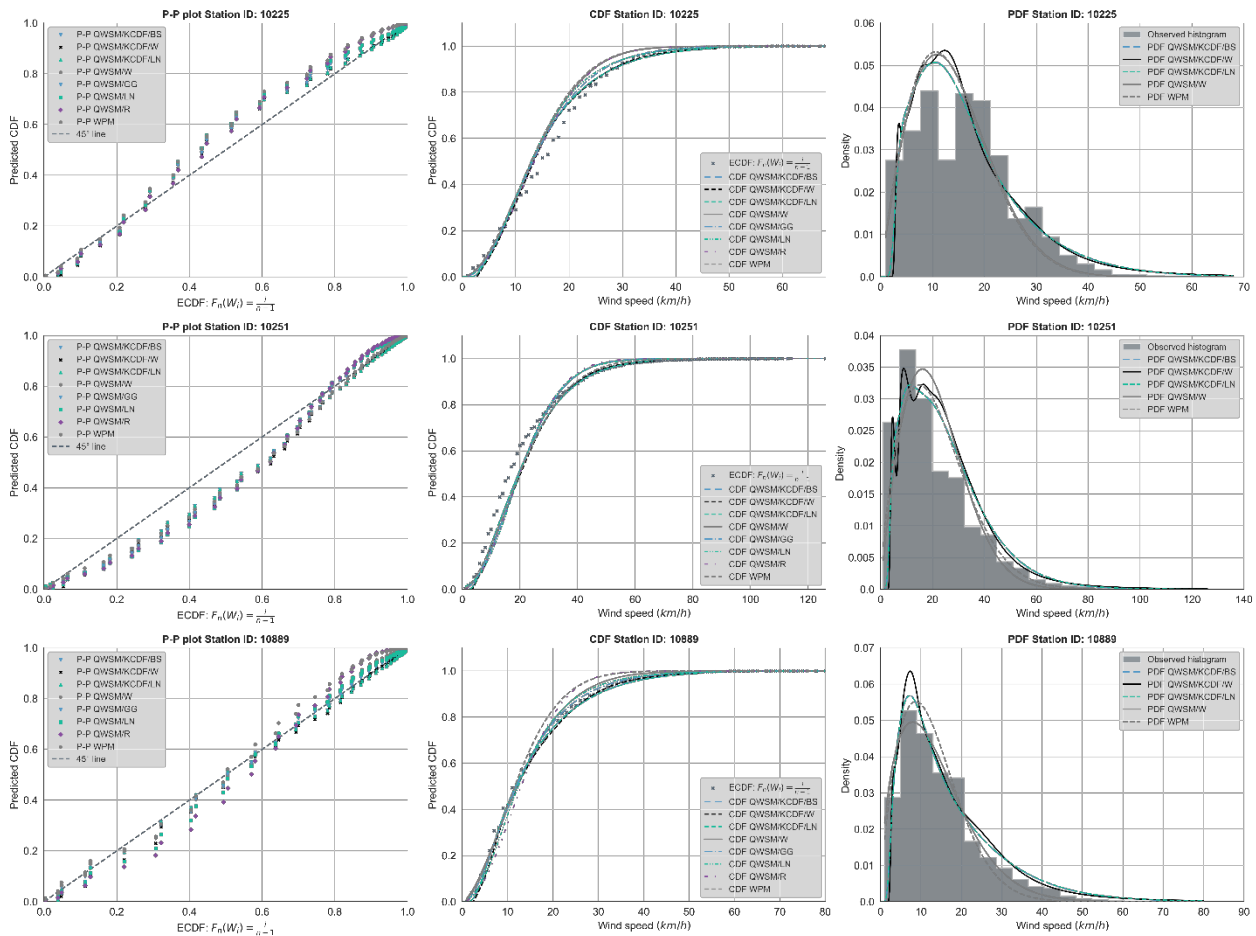
453 Figure 6: Performance metrics for observed WSQ and estimated WSQ using QWSM and WPM (validation
454 set)

455 In Figure 7, the P-P plot, the CDF, and the probability density function (PDF) plot of 3 validation samples
456 are presented for illustration purposes. These plots offer a comprehensive visual analysis of the actual
457 and estimated WS distribution agreement. Recall that QWSM/W was selected as the target distribution
458 to estimate the optimal bandwidth for all KCDF. However, it is observed that the kernel PDFs exhibited
459 more flexibility than QWSM/W. The W kernel demonstrated more flexibility than the BS and LN kernels,
460 while both gave an almost identical PDF.

461

462

463



465

466 Figure 7: PP plot, CDF plot and PDF plot of estimated wind speed probability distributions

467

468 5. Discussion

469 The comparison of the regression models indicates that the non-linear model (XGB) outperformed the
 470 linear model (LR) for the estimation of WSQ and the W parameters. The superior performance of the
 471 XGB model suggests that there are non-linear associations and interactions between the covariates and
 472 the WS response variables (WSQ and W parameters). The XGB model can effectively capture these non-
 473 linear relationships, leading to more accurate and precise estimates than the linear model. There is
 474 potential for further improvement in the performance of the XGB model by conducting a more
 475 comprehensive hyperparameter tuning. A random search was employed for the XGB hyperparameter

476 tuning and proved sufficient to demonstrate the superiority of the XGB model over the LR model.
477 However, a more extensive hyperparameter tuning process, such as grid search or Bayesian optimization
478 [52], could be conducted to thoroughly search for the optimal combination of hyperparameters that
479 maximizes the model's performance.

480 The study also found that MRMR-PC was more effective for FS than MRMR-MI. MI can assess linear and
481 nonlinear dependencies between variables, and it was initially expected that combining MRMR-MI with
482 XGB would outperform the combination of MRMR-PC with XGB. However, similar results were observed
483 by Ren, et al. [53] in the field of hydrology. The authors discovered that a FS method based on the partial
484 Pearson correlation outperformed FS methods based on MI (including MRMR-MI) when applied with
485 linear and nonlinear regression models for monthly streamflow forecasting. The study attributed these
486 results to the possibility that the relationship between the covariates and the target variable in their
487 models exhibited more linearity than nonlinearity. Similar conclusions may be formulated in this study,
488 suggesting that the gain in performance achieved using the XGB could also be attributed to other
489 characteristics of the models, such as its robustness against redundant features and collinearity within
490 the features set. Despite these findings, it is still recommended to evaluate different FS methods.
491 Different scenarios or datasets may yield different results.

492 It is well known that wind speed and other climatic variables like humidity, pressure, and temperature
493 are interconnected. The main challenge in using climatic variables for estimating wind speed at
494 unsampled locations is that those variables should also be unavailable. Gridded climate data can be used
495 as an alternative source of climatic covariates. This study only used gridded climate data of long-term
496 temperature trends as climatic covariates. Investigating the applicability of other gridded climate data as
497 covariates for WS distribution mapping in future studies is recommended.

498 Veronesi, et al. [8] reviewed the performance of physical and statistical methods for wind resources
499 assessment. They found that most studies applying statistical methods reported an RMSE of around 1
500 m/s on their validation set when considering the central tendency of the wind speed distribution (ex.:
501 mean). In the current study, the average RMSE for estimating the median wind speed obtained was 3.28
502 km/h (0.87 m/s), and the average MAE was 2.62 km/h (0.69 m/s). These results seem to agree with
503 previous studies. However, as was pointed out by Veronesi, et al. [8], results from different studies are
504 generally difficult to compare as different datasets, regions and techniques were covered in these
505 studies.

506 In general, based on the evaluation of the GOF, QWSM demonstrated a better fit compared to WPM.
507 This result may be explained by the fact that the estimation of the WS distribution from WSQ may be less
508 sensitive to mapping error compared to WPM. For instance, in the case of the WPM, minor errors in
509 mapping the W parameter could have disproportionate effects on the overall resulting shape of the wind
510 speed distribution. In contrast, with the QWSM, the implications of mapping errors are less severe, as
511 inaccuracies in wind speed quantile mapping seemed to have a smaller impact on the overall
512 distribution's shape. Consequently, the QWSM approach exhibits enhanced robustness against errors in
513 mapping, rendering it a more dependable framework for wind speed distribution mapping.

514 The non-parametric approach with the BS and LN KCDF gave slightly better results than the parametric
515 approach when considering the Kolmogorov-Smirnov statistic. The non-parametric method does not
516 require fixing a regional distribution and can adequately recover the WS distribution from the estimated
517 quantiles. Parametric methods require fitting the data to a specific probability distribution family, which
518 may introduce bias if the assumed distribution does not align with the underlying distribution. Another
519 potential source of bias common to both methods (i.e.: QWSM, WPM) is related to the regression
520 models used to estimate either the WSQ or the RD parameters. It should be noted that the bulk of the
521 bias of the QWSM + KCDF method arises from the regression model used to map the WSQ in the region.

522 Thus, the non-parametric approach can reduce potential biases by minimizing the assumptions. The
523 proposed approach becomes particularly interesting in regions where the wind regime exhibits
524 significant variations, and no single distribution family is suitable for all locations within the region. With
525 their constraints, parametric methods may struggle to capture the diversity of complex patterns that can
526 be present in such regions. In contrast, with its flexibility, the non-parametric approach can be more
527 appropriate and should yield more accurate results. Alternatively, it is possible to segregate the regions
528 into sub-regions and select a different RD for each sub-region. However, this would reduce the number
529 of samples used to learn the relationship between the covariates and the RD parameters, potentially
530 leading to a loss in performance. For WS values located in the distribution's tails (for instance, extreme
531 values), opting for the QWSM method with parametric distribution functions would be more suitable.
532 This recommendation is based on the finding that these parametric approaches exhibited superior
533 performance compared to non-parametric approaches in this case.

534 Mapping the WSQ in this study involved extracting the quantiles from the time series and then using a
535 regression model that estimates the conditional mean of the quantiles given the covariates. An
536 alternative approach could be directly estimating the conditional quantiles using a quantile regression
537 [54-57] model incorporating the covariates. Quantile regression is a statistical technique that allows
538 estimating specific quantiles of the response variable rather than focusing solely on the conditional
539 mean.

540 The main drawback of the QWSM approach is that the number of independent variables (quantiles) that
541 need to be mapped to recover the WS distribution would often be superior to the number of the RD
542 parameters that require mapping in the WPM approach. Fitting these individual regression models can
543 become time-consuming and resource intensive. However, some quantile regression models can
544 simultaneously estimate multiple quantiles [57, 58], providing a more efficient approach compared to
545 building separate regression models for each quantile. Also, when estimating multiple quantiles

546 simultaneously, additional constraints can be formulated to enforce monotonicity [59] and avoid the
547 issue of quantile crossing that arises when estimating the quantiles independently. It is worth
548 mentioning that a gradient-boosting model [60] was recently proposed to simultaneously estimate the
549 parameters of a probability distribution conditioned on some covariates. This model could be used to
550 estimate the parameters of a RD simultaneously rather than building an independent model for each
551 parameter.

552 Modern wind turbine hub heights vary between 80m and 100m, while wind speed data are
553 conventionally collected at 10m at meteorological stations. As a result, a technique for extrapolating
554 wind speed data to hub height becomes necessary (ex.: the power law). Such techniques can extend the
555 method proposed in this study to map wind speed distribution at hub height. Nevertheless, it is worth
556 noting that such extrapolation introduces a notable increment in the uncertainty of the outcomes.

557 Jung and Schindler [26] proposed a technique for mapping wind shear distribution, allowing the wind
558 speed distribution to be mapped at any standard hub height. Jung and Schindler [26] selected the Dagum
559 family distribution to represent the wind shear distribution. In future research, the non-parametric
560 approach proposed in this study could be adapted to map wind shear distribution without prior
561 assumptions about its distribution. Also, future studies can explore the possibility of extending the
562 proposed approach to other types of climatic variables, such as temperature and solar irradiation.

563 The approach proposed in this study can provide valuable information to estimate wind resources over a
564 large area during a prospecting phase. Once an area that meets the necessary socio-economic
565 requirements and showcases sufficient wind potential is identified, alternative methods are available to
566 evaluate the wind flow at the microscale. An example of such an approach involves conducting wind flow
567 simulations via Computational Fluid Dynamics (CFD), especially in complex terrain [61]. The
568 implementation of a CFD model requires the provision of initial wind data, which can be sourced from

569 outputs generated by Numerical Weather Prediction [NWP: 62, 63, 64]. NWP models entail considerable
570 computational costs compared to statistical methodologies proposed herein. A compelling avenue of
571 research would involve comparing the performance of NWP and statistical models for CFD model input
572 and developing methods to combine statistical and CFD models to assess microscale wind flow dynamics.

573

574 6. Conclusion

575 A fully non-parametric approach was developed to map wind speed distribution. The new method was
576 compared to a more traditional approach based on mapping the parameters of a regional distribution.
577 The results of the comparative analysis highlighted the superiority of the proposed approach. The main
578 conclusions of the paper are summarized as follow:

- 579 • The non-parametric approach is more practical as it does not require fitting and evaluating
580 several distribution functions to the available wind speed data. In the proposed method, wind
581 speed quantiles can be easily extracted from the time series and mapped using suitable
582 machine-learning techniques. At any location in the study area, the entire wind speed
583 distribution can be recovered from the estimated wind speed quantile by fitting asymmetric
584 kernel estimators. The proposed approach is free from any assumption on the wind speed
585 probability distribution family in the region that can bias the analysis. The non-parametric
586 approach is recommended for mapping wind speed distribution in regions with a highly variable
587 wind regime. The analysis indicates that the fully non-parametric approach improved the
588 Kolmogorov-Smirnov statistic by 9% on average during validation.
- 589 • Compared to the regional distribution parameter mapping approach, quantile-based wind speed
590 distribution mapping can be slower to implement as it requires the estimation of multiple wind
591 speed quantiles. However, with the advancement in quantile regression models, it is possible to

592 build a single regression model to predict multiple quantiles. This type of quantile regression
593 model should reduce the computational burden associated with the proposed approach.

- 594 • The Gradient boosting trees model outperformed the multilinear regression model for mapping
595 wind speed quantiles and the Weibull parameters. At the same time, feature selection based on
596 the Pearson correlation coefficient was more effective than the Mutual information. Utilizing the
597 Gradient Boosting Trees model and feature selection based on the Pearson correlation
598 coefficient resulted in a 23% improvement in R^2 during validation compared to the second-best
599 model for estimating wind speed quantiles.
- 600 • It should be noted that symmetric kernels could also be fitted to the estimated wind speed
601 quantiles, with some probabilities associated to small negative wind speed values. Using an
602 asymmetric kernel effectively avoids probability leakage at the boundary of the lower tail of the
603 wind speed probability distribution.
- 604 • The proposed approach is easily portable to regions with sparsely available wind speed
605 measuring stations. The other data sources used in the study (ex.: DEM and land use map) are
606 often freely accessible from global datasets covering most regions of the world.

607 7. Acknowledgments

608 The authors thank the Natural Sciences and Engineering Research Council of Canada (NSERC) and the
609 Canada Research Chair Program for funding this research. The wind speed data was freely accessed from
610 Environment and Climate Change Canada database, and the Japan Aerospace Exploration Agency provided
611 free access to the digital elevation model. The authors would like to extend their gratitude to the Editor,
612 Prof. Nesreen Ghaddar, and three anonymous reviewers for their comments and suggestions, which
613 significantly improved the quality of the paper.

614

615 Appendix I Statistics of the estimated wind speed quantiles

616 Table 6: Statistics of the estimated wind speed quantiles

Percentile	Mean (km/h)	Std (km/h)	Min (km/h)	25% (km/h)	50% (km/h)	75% (km/h)	Max (km/h)
5	4.2	1.6	1	3	4	5	9
12.5	6.4	2.2	2	5	6	7	13
20	8.2	2.9	3	6	7	9.5	18
27.5	9.9	3.4	4	7	9	12	20
35	11.6	4.0	4	9	11	14.5	24
42.5	13.4	4.4	5	11	13	16.5	28
50	15.2	4.9	6	12	15	19	31
57.5	17.1	5.4	6	13	17	20	35
65	19.4	6.1	7	15	19	23.5	39
72.5	21.9	6.7	7	17	21	26	44
80	24.8	7.6	9	19	24	30	51
87.5	29.0	8.7	11	22.5	28	35	59
95	36.2	10.9	14	28.5	35	44	74

617

618 Appendix II. Wind speed covariates

619 Table 7: Overview of the WS covariates

Predictor	Description	Spatial scale
Altitude	Altitude of the location in meter.	
Aspect	Slope orientation in degree.	100m, 500m, 1000m, 1500m, 2000m
Deviation from mean elevation	Difference between the grid cell elevation and the mean of its neighbouring cells normalized by the standard deviation.	100m, 500m, 1000m, 1500m, 2000m
Difference from cell mean elevation	Difference between the grid cell elevation and the mean of its neighbouring cells.	100m, 500m, 1000m, 1500m, 2000m
Difference of Gaussian	Difference between two copies of the DEM smoothed with two different gaussian kernel. Measure land surface curvature.	(100m, 500m), (100m, 1000m), (300m, 500m), (1000m, 2000m), (100m, 2000m), (500m, 1000m), (1000m, 1500m), (500m, 2000m)
Distance to coast	The location distance to the coast	
Elevation percentile	Percentile of the grid cell elevation relative to the neighbouring cells.	100m, 500m, 1000m, 1500m, 2000m

Gaussian curvature	Product between the maximal and the minimal curvature. Measure of surface curvature [65].	100m, 500m, 1000m, 1500m, 2000m
Geographical coordinates	Geographical coordinates of the location.	
geomorphologic phenotypes (geomorphons)	Landform element classification with the geomorphons-based method [66].	
Laplacian of Gaussian	Derivative filter used to highlight location of rapid elevation change.	100m, 500m, 1000m, 1500m, 2000m
Maximal curvature	Measure of surface curvature [67].	100m, 500m, 1000m, 1500m, 2000m
Mean curvature	Measure of surface curvature [67].	100m, 500m, 1000m, 1500m, 2000m
minimal curvature	Measure of surface curvature [65].	100m, 500m, 1000m, 1500m, 2000m
Pennock landform class	Landform classification based on the slope and curvature of the grid cell [68].	
plan curvature	Measure of surface curvature [65].	100m, 500m, 1000m, 1500m, 2000m
Relative topographical position	Normalized measure of the grid cell elevation relative to its neighbouring cells.	100m, 500m, 1000m, 1500m, 2000m
Ruggedness index	A measure of the local terrain heterogeneity [66, 69]	100m, 500m, 1000m, 1500m, 2000m
Slope	Slope at the grid cell.	100m, 500m, 1000m, 1500m, 2000m
Standard deviation of slope	Measure of surface roughness [70].	100m, 500m, 1000m, 1500m, 2000m
Surface area ratio	Measure of the surface roughness [71].	100m, 500m, 1000m, 1500m, 2000m
Surface roughness length	Surface roughness length estimated from land use map.	100m, 500m, 1000m, 1500m, 2000m
tangential curvature	Measure of surface curvature [65].	100m, 500m, 1000m, 1500m, 2000m
Total curvature	Measure of surface curvature.	100m, 500m, 1000m, 1500m, 2000m
Temperature trend	Seasonal and annual trends of mean temperature change between 1948-2018.	

620

621

622 **References**

- 623 [1] Y. Zhou, P. Luckow, S.J. Smith, L. Clarke. Evaluation of Global Onshore Wind Energy Potential and
624 Generation Costs. *Environmental Science & Technology*. 46 (2012) 7857-64. [10.1021/es204706m](https://doi.org/10.1021/es204706m)
- 625 [2] Global Wind Energy Council. GWEC Global Wind Report 2022. GWEC Global Wind Report. Global
626 Wind Energy Council, Brussels, Belgium, 2022.
- 627 [3] C. Jung, D. Schindler, J. Laible. National and global wind resource assessment under six wind turbine
628 installation scenarios. *Energy Conversion and Management*. 156 (2018) 403-15.
629 <https://doi.org/10.1016/j.enconman.2017.11.059>
- 630 [4] F. Houndekindo, T.B.M.J. Ouarda. Statistical approaches for wind speed estimation at ungauged or
631 partially gauged locations, review, and open questions (Under review). Institut national de la recherche
632 scientifique, Centre Eau Terre Environnement. (2023).
- 633 [5] W. Luo, M.C. Taylor, S.R. Parker. A comparison of spatial interpolation methods to estimate
634 continuous wind speed surfaces using irregularly distributed data from England and Wales. *International
635 Journal of Climatology*. 28 (2008) 947-59. <https://doi.org/10.1002/joc.1583>
- 636 [6] W. Ye, H.P. Hong, J.F. Wang. Comparison of Spatial Interpolation Methods for Extreme Wind Speeds
637 over Canada. *Journal of Computing in Civil Engineering*. 29 (2015) 04014095.
638 [doi:10.1061/\(ASCE\)CP.1943-5487.0000429](https://doi.org/10.1061/(ASCE)CP.1943-5487.0000429)
- 639 [7] V. Nelson, K.r.e. Starcher. *Wind Energy: Renewable Energy and the Environment*. CRC Press, Boca
640 raton, Floride USA, 2018.
- 641 [8] F. Veronesi, S. Grassi, M. Raubal. Statistical learning approach for wind resource assessment.
642 *Renewable and Sustainable Energy Reviews*. 56 (2016) 836-50.
643 <https://doi.org/10.1016/j.rser.2015.11.099>

- 644 [9] C. Jung, D. Schindler. Wind speed distribution selection – A review of recent development and
645 progress. *Renewable and Sustainable Energy Reviews*. 114 (2019) 109290.
646 <https://doi.org/10.1016/j.rser.2019.109290>
- 647 [10] O. Tsvetkova, T.B.M.J. Ouarda. Use of the Halphen distribution family for mean wind speed
648 estimation with application to Eastern Canada. *Energy Conversion and Management*. 276 (2023) 116502.
649 [10.1016/j.enconman.2022.116502](https://doi.org/10.1016/j.enconman.2022.116502)
- 650 [11] C. Jung. High Spatial Resolution Simulation of Annual Wind Energy Yield Using Near-Surface Wind
651 Speed Time Series. *Energies*. 9 (2016) 344. doi:10.3390/en9050344
- 652 [12] M. Laib, M. Kanevski. Spatial Modelling of Extreme Wind Speed Distributions in Switzerland. *Energy*
653 *Procedia*. 97 (2016) 100-7. <https://doi.org/10.1016/j.egypro.2016.10.029>
- 654 [13] T.B.M.J. Ouarda, C. Charron. On the mixture of wind speed distribution in a Nordic region. *Energy*
655 *Conversion and Management*. 174 (2018) 33-44. <https://doi.org/10.1016/j.enconman.2018.08.007>
- 656 [14] J. Zhou, E. Erdem, G. Li, J. Shi. Comprehensive evaluation of wind speed distribution models: A case
657 study for North Dakota sites. *Energy Conversion and Management*. 51 (2010) 1449-58.
658 <https://doi.org/10.1016/j.enconman.2010.01.020>
- 659 [15] B. Safari. Modeling wind speed and wind power distributions in Rwanda. *Renewable and Sustainable*
660 *Energy Reviews*. 15 (2011) 925-35. <https://doi.org/10.1016/j.rser.2010.11.001>
- 661 [16] N. Aries, S.M. Boudia, H. Ounis. Deep assessment of wind speed distribution models: A case study of
662 four sites in Algeria. *Energy Conversion and Management*. 155 (2018) 78-90.
663 <https://doi.org/10.1016/j.enconman.2017.10.082>
- 664 [17] O. Alavi, K. Mohammadi, A. Mostafaeipour. Evaluating the suitability of wind speed probability
665 distribution models: A case of study of east and southeast parts of Iran. *Energy Conversion and*
666 *Management*. 119 (2016) 101-8. <https://doi.org/10.1016/j.enconman.2016.04.039>

667 [18] T.B.M.J. Ouarda, C. Charron, J.Y. Shin, P.R. Marpu, A.H. Al-Mandoos, M.H. Al-Tamimi, et al.
668 Probability distributions of wind speed in the UAE. *Energy Conversion and Management*. 93 (2015) 414-
669 34. <http://doi.org/10.1016/j.enconman.2015.01.036>

670 [19] Q. Han, S. Ma, T. Wang, F. Chu. Kernel density estimation model for wind speed probability
671 distribution with applicability to wind energy assessment in China. *Renewable and Sustainable Energy*
672 *Reviews*. 115 (2019) 109387. <https://doi.org/10.1016/j.rser.2019.109387>

673 [20] S. Węglarczyk. Kernel density estimation and its application. *ITM Web Conf*. 23 (2018).

674 [21] H.A. Mombeni, B. Mansouri, M. Akhoond. Asymmetric kernels for boundary modification in
675 distribution function estimation. *REVSTAT-Statistical Journal*. 19 (2021) 463–84–84.

676 [22] M. Hirukawa. *Asymmetric Kernel Smoothing: Theory and Applications in Economics and Finance*.
677 Springer Nature Singapore, Singapore, 2018.

678 [23] R.J. Hyndman, Y. Fan. Sample quantiles in statistical packages. *The American Statistician*. 50 (1996)
679 361-5.

680 [24] R.D. Reiss. *Approximate Distributions of Order Statistics: With Applications to Nonparametric*
681 *Statistics*. Springer New York 1989.

682 [25] Y. He, R. Liu, H. Li, S. Wang, X. Lu. Short-term power load probability density forecasting method
683 using kernel-based support vector quantile regression and Copula theory. *Applied Energy*. 185 (2017)
684 254-66. <https://doi.org/10.1016/j.apenergy.2016.10.079>

685 [26] C. Jung, D. Schindler. 3D statistical mapping of Germany's wind resource using WSWS. *Energy*
686 *Conversion and Management*. 159 (2018) 96-108. <https://doi.org/10.1016/j.enconman.2017.12.095>

687 [27] T. B.M.J. Ouarda, C. Charron, A. St-Hilaire. Regional estimation of river water temperature at
688 ungauged locations. *Journal of Hydrology X*. (2022). 10.1016/j.hydroa.2022.100133

689 [28] J.H. Friedman. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of*
690 *Statistics*. 29 (2001) 1189-232.

691 [29] C. Ding, P. Hanchuan. Minimum redundancy feature selection from microarray gene expression
692 data. *J Bioinform Comput Biol.* 3 (2005) 185-205. [10.1142/s0219720005001004](https://doi.org/10.1142/s0219720005001004)

693 [30] F. Houndekindo, T.B.M.J. Ouarda. Comparative study of feature selection methods for wind speed
694 estimation at ungauged locations. *Energy Conversion and Management.* 291 (2023) 117324.
695 [10.1016/j.enconman.2023.117324](https://doi.org/10.1016/j.enconman.2023.117324)

696 [31] I. Guyon, A. Elisseeff. An introduction to variable and feature selection. *Journal of machine learning*
697 *research.* 3 (2003) 1157-82.

698 [32] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, et al. Scikit-learn: Machine
699 learning in Python. *the Journal of machine Learning research.* 12 (2011) 2825-30.

700 [33] L. Grinsztajn, E. Oyallon, G. Varoquaux. Why do tree-based models still outperform deep learning on
701 typical tabular data? , *Thirty-sixth Conference on Neural Information Processing Systems Datasets and*
702 *Benchmarks Track2022.*

703 [34] T. Hastie, R. Tibshirani, J.H. Friedman, J.H. Friedman. *The elements of statistical learning: data*
704 *mining, inference, and prediction.* Springer, New York, 2009.

705 [35] L. Ye, B. Dai, Z. Li, M. Pei, Y. Zhao, P. Lu. An ensemble method for short-term wind power prediction
706 considering error correction strategy. *Applied Energy.* 322 (2022) 119475.
707 <https://doi.org/10.1016/j.apenergy.2022.119475>

708 [36] Y. Sun, D. Zhu, Y. Li, R. Wang, R. Ma. Spatial modelling the location choice of large-scale solar
709 photovoltaic power plants: Application of interpretable machine learning techniques and the national
710 inventory. *Energy Conversion and Management.* 289 (2023) 117198.
711 <https://doi.org/10.1016/j.enconman.2023.117198>

712 [37] J. Lee, W. Wang, F. Harrou, Y. Sun. Reliable solar irradiance prediction using ensemble learning-
713 based models: A comparative study. *Energy Conversion and Management.* 208 (2020) 112582.
714 <https://doi.org/10.1016/j.enconman.2020.112582>

- 715 [38] T. Chen, C. Guestrin. XGBoost: A Scalable Tree Boosting System. Proceedings of the 22nd ACM
716 SIGKDD International Conference on Knowledge Discovery and Data Mining. Association for Computing
717 Machinery, San Francisco, California, USA, 2016. pp. 785–94.
- 718 [39] R. Turner, D. Eriksson, M. McCourt, J. Kiili, E. Laaksonen, Z. Xu, I. Guyon. Bayesian Optimization is
719 Superior to Random Search for Machine Learning Hyperparameter Tuning: Analysis of the Black-Box
720 Optimization Challenge 2020. in: E. Hugo Jair, H. Katja, (Eds.), Proceedings of the NeurIPS 2020
721 Competition and Demonstration Track. PMLR, Proceedings of Machine Learning Research, 2021. pp. 3--
722 26.
- 723 [40] J. Bergstra, Y. Bengio. Random search for hyper-parameter optimization. Journal of machine
724 learning research. 13 (2012).
- 725 [41] T.B.M.J. Ouarda, C. Charron. Non-stationary statistical modelling of wind speed: A case study in
726 eastern Canada. Energy Conversion and Management. 236 (2021) 114028.
727 <https://doi.org/10.1016/j.enconman.2021.114028>
- 728 [42] P. Lafaye de Micheaux, F. Ouimet. A Study of Seven Asymmetric Kernels for the Estimation of
729 Cumulative Distribution Functions. Mathematics2021.
- 730 [43] T.B.M.J. Ouarda, C. Charron, F. Chebana. Review of criteria for the selection of probability
731 distributions for wind speed data and introduction of the moment and L-moment ratio diagram
732 methods, with a case study. Energy Conversion and Management. 124 (2016) 247-65.
733 <http://dx.doi.org/10.1016/j.enconman.2016.07.012>
- 734 [44] M.B. Wilk, R. Gnanadesikan. Probability plotting methods for the analysis for the analysis of data.
735 Biometrika. 55 (1968) 1-17. 10.1093/biomet/55.1.1
- 736 [45] R. Horst. The Weibull Distribution: A Handbook. Chapman and Hall/CRC, New York, 2008.

737 [46] F.G. Akgül, B. Şenoğlu, T. Arslan. An alternative distribution to Weibull for modeling the wind speed
738 data: Inverse Weibull distribution. *Energy Conversion and Management*. 114 (2016) 234-40.
739 <https://doi.org/10.1016/j.enconman.2016.02.026>

740 [47] E.C. Morgan, M. Lackner, R.M. Vogel, L.G. Baise. Probability distributions for offshore wind speeds.
741 *Energy Conversion and Management*. 52 (2011) 15-26. <https://doi.org/10.1016/j.enconman.2010.06.015>

742 [48] J.B. Lindsay. The Whitebox Geospatial Analysis Tools project and open-access GIS. GIS Research UK
743 22nd Annual Conference. The University of Glasgow, University of Glasgow, 2014.

744 [49] T. Tadono, H. Ishida, F. Oda, S. Naito, K. Minakawa, H. Iwamoto. Precise Global DEM Generation by
745 ALOS PRISM. *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*. II4
746 (2014) 71-6. [10.5194/isprsannals-II-4-71-2014](https://doi.org/10.5194/isprsannals-II-4-71-2014)

747 [50] R. Latifovic, D. Pouliot, I. Olthof. Circa 2010 Land Cover of Canada: Local Optimization Methodology
748 and Product Development. *Remote Sensing*. 9 (2017) 1098.

749 [51] J. Wiernga. Representative roughness parameters for homogeneous terrain. *Boundary-Layer*
750 *Meteorology*. 63 (1993) 323-63. [10.1007/BF00705357](https://doi.org/10.1007/BF00705357)

751 [52] J. Wu, X.-Y. Chen, H. Zhang, L.-D. Xiong, H. Lei, S.-H. Deng. Hyperparameter Optimization for
752 Machine Learning Models Based on Bayesian Optimization. *Journal of Electronic Science and*
753 *Technology*. 17 (2019) 26-40. <https://doi.org/10.11989/JEST.1674-862X.80904120>

754 [53] K. Ren, W. Fang, J. Qu, X. Zhang, X. Shi. Comparison of eight filter-based feature selection methods
755 for monthly streamflow forecasting – Three case studies on CAMELS data sets. *Journal of Hydrology*. 586
756 (2020) 124897. <https://doi.org/10.1016/j.jhydrol.2020.124897>

757 [54] R. Koenker. Quantile Regression: 40 Years On. *Annual Review of Economics*. 9 (2017) 155-76.
758 [10.1146/annurev-economics-063016-103651](https://doi.org/10.1146/annurev-economics-063016-103651)

759 [55] B. Nasri, T. Bouezmarni, A. St-Hilaire, T.B.M.J. Ouarda. Non-stationary hydrologic frequency analysis
760 using B-spline quantile regression. *J Hydrol.* 554 (2017) 532-44.
761 <https://doi.org/10.1016/j.jhydrol.2017.09.035>

762 [56] D. Ouali, F. Chebana, T. Ouarda. Quantile Regression in Regional Frequency Analysis: A Better
763 Exploitation of the Available Information. *Journal of Hydrometeorology.* 17 (2016). 10.1175/JHM-D-15-
764 0187.1

765 [57] N. Meinshausen, G. Ridgeway. Quantile regression forests. *Journal of machine learning research.* 7
766 (2006).

767 [58] Y. Liu, Y. Wu. Simultaneous multiple non-crossing quantile regression estimation using kernel
768 constraints. *Journal of Nonparametric Statistics.* 23 (2011) 415-37. 10.1080/10485252.2010.537336

769 [59] A.J. Cannon. Non-crossing nonlinear regression quantiles by monotone composite quantile
770 regression neural network, with application to rainfall extremes. *Stochastic Environmental Research and*
771 *Risk Assessment.* 32 (2018) 3207-25. 10.1007/s00477-018-1573-6

772 [60] T. Duan, A. Anand, D.Y. Ding, K.K. Thai, S. Basu, A. Ng, A. Schuler. Ngboost: Natural gradient boosting
773 for probabilistic prediction. *International Conference on Machine Learning.* PMLR2020. pp. 2690-700.

774 [61] X.-Y. Tang, S. Zhao, B. Fan, J. Peinke, B. Stoevesandt. Micro-scale wind resource assessment in
775 complex terrain based on CFD coupled measurement from multiple masts. *Applied Energy.* 238 (2019)
776 806-15. <https://doi.org/10.1016/j.apenergy.2019.01.129>

777 [62] P. Beaucage, M.C. Brower, J. Tensen. Evaluation of four numerical wind flow models for wind
778 resource mapping. *Wind Energy.* 17 (2014) 197-208. <https://doi.org/10.1002/we.1568>

779 [63] R.E. Keck, N. Sondell. Validation of uncertainty reduction by using multiple transfer locations for
780 WRF–CFD coupling in numerical wind energy assessments. *Wind Energ Sci.* 5 (2020) 997-1005.
781 10.5194/wes-5-997-2020

782 [64] T. Simões, A. Estanqueiro. A new methodology for urban wind resource assessment. *Renewable*
783 *Energy*. 89 (2016) 598-605. <https://doi.org/10.1016/j.renene.2015.12.008>

784 [65] I.V. Florinsky. An illustrated introduction to general geomorphometry. *Progress in Physical*
785 *Geography: Earth and Environment*. 41 (2017) 723-52. 10.1177/0309133317733667

786 [66] J. Jasiewicz, T.F. Stepinski. Geomorphons — a pattern recognition approach to classification and
787 mapping of landforms. *Geomorphology*. 182 (2013) 147-56.
788 <https://doi.org/10.1016/j.geomorph.2012.11.005>

789 [67] J.P. Wilson. *Environmental applications of digital terrain modeling*. John Wiley & Sons 2018.

790 [68] D.J. Pennock, B.J. Zebarth, E. De Jong. Landform classification and soil distribution in hummocky
791 terrain, Saskatchewan, Canada. *Geoderma*. 40 (1987) 297-315. [https://doi.org/10.1016/0016-](https://doi.org/10.1016/0016-7061(87)90040-1)
792 [7061\(87\)90040-1](https://doi.org/10.1016/0016-7061(87)90040-1)

793 [69] S.J. Riley, S.D. DeGloria, R. Elliot. Index that quantifies topographic heterogeneity. *Intermountain*
794 *Journal of Sciences*. 5 (1999) 23-7.

795 [70] C.H. Grohmann, M.J. Smith, C. Riccomini. Multiscale Analysis of Topographic Surface Roughness in
796 the Midland Valley, Scotland. *IEEE Transactions on Geoscience and Remote Sensing*. 49 (2011) 1200-13.
797 10.1109/TGRS.2010.2053546

798 [71] J. Jenness. Calculating Landscape Surface Area from Digital Elevation Models. *Wildlife Society*
799 *Bulletin*. 32 (2004) 829-39. 10.2193/0091-7648(2004)032[0829:CLSAFD]2.0.CO;2

800