

Article

Not peer-reviewed version

A Hybrid Linear–Gaussian Process Framework with Adaptive Covariance Selection for Spatio-Temporal Wind Speed Forecasting

[Thinawanga Hangwani Tshisikhawe](#) , [Caston Sigauke](#) * , [Timotheus Brian Darikwa](#) , [Saralees Nadarajah](#)

Posted Date: 3 March 2026

doi: 10.20944/preprints202603.0262.v1

Keywords: cluster validation; Gaussian processes; hybrid models; spatio-temporal modelling; wind speed forecasting



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

A Hybrid Linear–Gaussian Process Framework with Adaptive Covariance Selection for Spatio-Temporal Wind Speed Forecasting

Thinawanga Hangwani Tshisikhawe ¹, Caston Sigauke ^{1,*}, Timotheous Brian Darikwa ²
and Saralees Nadarajah ³

¹ Department of Mathematical and Computational Sciences, University of Venda, South Africa

² Department of Statistics and Operations Research, University of Limpopo, South Africa

³ Department of Mathematics, Manchester University, United Kingdom

* Correspondence: caston.sigauke@univen.ac.za

Abstract

Accurate wind speed forecasting is critical for renewable energy planning and meteorological studies. However, wind behaviour is complex due to the presence of synoptic systems, terrain effects, and turbulence. This paper proposes a new model that uses a linear regression mean model and a Gaussian process for residual modelling. The monitoring stations were clustered by geographic coordinates and elevation. The Hopkins statistic was employed for cluster validation, while silhouette values were employed for cluster quality validation. It was found that for stations at high elevation located in the interior (Cluster 2), the GP model for residual modelling consistently improved wind forecast accuracy by up to 16.3%. However, for coastal stations at low elevation (Cluster 1), the GP model was not effective for residual modelling. This proves that the accuracy of GP residual modelling depends to a large extent on the wind regime.

Keywords: cluster validation; Gaussian processes; hybrid models; spatio-temporal modelling; wind speed forecasting

1. Introduction

1.1. Overview

Wind speed forecasting techniques are applied in many fields that involve integrating green energy, power systems, and climate change analysis. However, wind exhibits complex, intermittent spatial and temporal dynamics that depend on atmospheric conditions, topography, and local turbulent eddies. While linear models effectively explain global behaviour and ignore any local dependencies, data-driven forecasting techniques lack interpretability and any quantification of uncertainty.

In the context of spatio-temporal forecasting, Gaussian Process (GP) approaches offer the advantage of a probabilistic treatment of uncertainty while enabling the covariance structure to be modelled economically. Despite these advantages, GPs tend to be computationally expensive, especially for large datasets, and are also computationally redundant. A hybrid approach to combining GPs is to tailor the models to the residuals of a linear mean model. Generally, wind speed is characterised differently across the country, with trends ranging from low-elevation coastal regions to high-elevation interior regions. Therefore, unless the spatial variability is modelled appropriately, the forecasting could be inaccurate. To address the above challenges, a model is proposed that uses a hybrid method to model regional wind speed. The model partitions regional wind speeds into a mean model and a GP model for the errors.

1.2. Literature Review

Accurate forecasting of wind speeds is critical for integrating wind energy into the grid and for various meteorological needs. Unfortunately, the chaotic behaviour of wind, influenced by large-scale atmospheric conditions, local topography, and turbulence, makes wind speed forecasting a complex problem. This review traces the chronological and thematic developments of wind speed forecasting techniques, from earlier physical and statistical models to modern hybrid and probabilistic models. It focuses on the developments of Gaussian Process Regression (GPR) and the utilisation of clustering techniques, both of which are significant to the proposed research. A brief overview of the major models is presented in Table 1.

1.2.1. Classical Methods: Physical and Statistical Models

Traditionally, wind speed forecasting was carried out using Numerical Weather Prediction (NWP) models, which simulate atmospheric physics to predict wind speeds. These are certainly essential in wind speed forecasting; however, these models are subject to systematic errors and lack spatial resolution, which are often overcome by post-processing these results to increase their accuracy locally [1]. To overcome these problems, statistical time-series models such as ARIMA and SARIMA were widely applied. These models are effective in capturing linear trends and seasonality in wind speeds; however, they are not effective in capturing rapid, non-linear changes in wind speeds, which are characteristic of wind speeds in general [2]. More adaptive models, such as Kalman Filters, were then applied to correct errors in short-term forecasts; however, these filters rely on linear assumptions and are prone to noise, limiting their application to wind speed forecasting over longer time scales and in spatially heterogeneous domains [3].

1.2.2. The Ascent of Machine Learning and Deep Learning

The limitations of linear models triggered the development and application of machine learning (ML) and deep learning (DL) techniques. ML and DL models, including random forests, Convolutional Neural Networks (CNNs), and Long Short-Term Memory (LSTM), are very effective at learning complex, non-linear relationships directly from historical data, often resulting in significant accuracy gains in forecasts [4]. The hybrid model, combining NWP model results and LSTM, is also very effective for forecasting complex terrain, including non-linear relationships and temporal dependencies [5]. However, ML/DL-based models are computationally intensive, require careful tuning, and often struggle to generalise across different locations. Additionally, ML/DL models are considered "black box" models, which limits their interpretability, making interpretability a driving need for finding models that are effective, yet transparent and physically plausible [4,5].

1.2.3. Probabilistic Forecasting and Gaussian Process Regression

As the need to quantify uncertainty in grid operations increases, probabilistic forecasting has become a crucial area of study. Unlike point forecasts, probabilistic forecasting can offer predictive distributions for better risk-informed decision-making [6]. Bayesian methods, including Bayesian Neural Networks and Bayesian Optimised models, inherently offer uncertainty quantification. However, their site-specific calibration can be a hindrance for real-time applications [7].

Within this paradigm, Gaussian Process Regression (GPR) has emerged as a prominent non-parametric tool for probabilistic forecasting. GPR is particularly prized for its capacity to yield a mean estimate and a well-founded interval estimate, which capture the essential randomness of wind phenomena [8]. The efficacy of GPR is critically dependent on its covariance function (kernel), which embodies assumptions regarding the underlying process, e.g., smoothness, periodicity, etc. GPR has been successfully utilised in hybrid frameworks, e.g., coupling GPR with support vector regression to post-process NWP forecasts for short-term forecasting (1 to 6 hours ahead) [9]. Another hybrid model, which is a focus of this research, is a linear mean model coupled with a GPR model for forecasting residuals. This model balances interpretability with the power of GPR for modelling

complex, spatially correlated error structures. However, standard GPR is known to scale poorly to large datasets, requiring a cubic time complexity ($\mathcal{O}(n^3)$) for its implementation [9,10].

1.2.4. Spatio-Temporal Modelling and the Role of Clustering

The spatial correlations in wind fields suggest that using data from adjacent areas can improve local predictions. A more advanced version of GPR, known as kriging, which is analogous to geostatistics, has been employed to interpolate wind fields while accounting for spatial correlations. It provides predictions along with their uncertainties [9]. More recent models seek to address complex spatio-temporal interactions. For example, dynamic versions of GPR that incorporate local clustering for improving accuracy have been presented [11]. Bayesian models for spatio-temporal forecasting for renewable energy forecasting have also been presented, but their success relies on the availability of dense data [12].

A notable result of the recent research is the recognition that wind speed statistical properties are not necessarily homogeneous over vast geographical areas. Various regimes, depending on the presence or absence of factors such as coastal areas and altitude, exhibit unique properties, including levels of non-stationarity, turbulence, and autocorrelation. This means that rather than seeking an optimal model applicable over the global domain, the regime-based approach of dividing the domain into homogeneous areas is more advisable. Grouping stations by their physical properties, such as geographic coordinates and altitude, is a statistically sound approach for dividing the domain.

Recent research emphasises the potential of using the latest GP variants. For example, the application of hybrid GP approaches, which integrate spatial correlation with corrected NWP forecast information, has been found to offer better accuracy than the standard GP model [13]. Another GP model variant, the multi-task GP, is also found to be beneficial, particularly when applied to data-scarce locations, where information from data-rich locations can be exploited [14]. However, the current GP model variants, including the complex residual model, have the limitation that their added value is assumed to be applicable uniformly across the study area.

1.3. Summary of Literature Review

This research addresses the limitation of the complex residual model directly by proposing a hybrid linear-GP model, but not over the entire study area, but over geographically and elevation-based clusters. This is to evaluate the added value of the complex residual model across regimes, providing a physically informed, computationally efficient tool to improve the accuracy of wind speed forecasts.

A summary of the comparison of wind speed forecasting methods, including the key studies discussed, highlighting their temporal focus, uncertainty quantification, computational cost and limitations, is given in Table 1.

Table 1. Comparison of wind speed forecasting methods.

Method	Temporal Focus	Spatial Awareness	Uncertainty Quantification	Computational Cost	Key Limitations
NWP Models [1]	Medium-long range	Low	No	High	Systematic errors, limited resolution
SARIMA [2]	Short term	None	No	Low	Poor with non-linear data
Kalman Filters [3]	Short term	Low	Limited	Low-medium	Linearity assumptions
ML/DL [4,5]	Short-medium	None-medium	No	High	Black-box, site-specific, data-hungry
Bayesian Models [6,7]	Short-medium	None-medium	Yes	High-very high	Computationally intensive, calibration needed
GPR [8,9]	Short term	Medium	Yes	Very high ($O(n^3)$)	Scalability issues
Hybrid Linear-GPR [9,10]	Short term	Medium	Yes	Medium-high	Assumes uniform spatial benefit
Spatio-temporal GPR/Kriging [9,11]	Short term	High	Yes	Very high	Data density requirements
Cluster-based Hybrid (Proposed)	Short term	High	Yes	Medium	Regime identification, clustering sensitivity

1.4. Contribution and Research Highlights

Spatiotemporal wind forecasting models have a critical limitation in accounting for wake effects from wind turbines, which can significantly reduce forecast accuracy and the operational efficiency of wind farms [15]. Deep learning-based wind forecasting models have demonstrated strong capability to learn nonlinear temporal dependencies and improve short-term prediction accuracy; however, their black-box nature, high data and computational demands, and lack of inherent uncertainty quantification limit their interpretability and reliability for operational wind energy decision-making [16]. Additionally, the scalability of Gaussian Process Regression (GPR) limits its direct application in large wind farms, as computational demands reduce operational efficiency.

Multi-kernel Learning/Regression: Unified Spatio-Temporal Kernel which combines Matérn kernel (space) + Periodic kernel (time) to model the interactions between the turbines. This kernel identifies where and when interactions between the turbines occur. Bayesian uncertainty provides a measure of predictive variance on how confident a model is in its predictions. For grid operators, this helps quantify risk in power demand/supply forecasts, make data-driven decisions under uncertainty, prioritise responses to high-variance (less certain) scenarios, and improve the grid's resilience and reliability. In essence, it informs operators how much trust they can place in a forecast. GPR's superiority in uncertainty-aware forecasting is well-established, but scalability remains a hurdle. Bayesian methods (such as GPR, BNNs) are gaining traction over deterministic ML in energy applications.

The hybrid linear-Gaussian Process model presented here addresses key limitations of existing wind forecasting approaches by decomposing wind prediction into a deterministic linear component and a spatial GP residual component. This formulation preserves interpretability, explicitly models spatial dependence, and provides rigorous uncertainty quantification, making it particularly suitable for sparse wind station networks. Clustering stations based on elevation and geographic coordinates (longitude and latitude) further improves model performance by capturing location-specific wind patterns and local topographic effects. Model accuracy is enhanced through careful kernel selection, where candidate spatial and temporal kernels are evaluated, and hyperparameter optimisation via

grid search ensures that the Gaussian Process component is well-calibrated to the underlying data. Together, these strategies improve predictive performance and provide reliable uncertainty estimates.

The rest of the paper is organised as follows. Section 2 briefly discusses materials and methods, including the Hopkins variable selection method, temporal GP regression, linear temporal regression, the hybrid, and the evaluation metrics. Section 3 presents the results, and the discussion is presented in Section 4, and the conclusion in Section 5.

2. Materials and Methods

2.1. Data Description and Clustering

The data used in this study are from the Wind Atlas South Africa (WASA) at <https://wasadata.csir.co.za/wasa1/WASAData>. The WASA data was obtained from the website [17]. Wind speed data from 10 stations were collected, each with multiple meteorological predictors (e.g., temperature, pressure). The variables included are as follows: Wind speed (m/s) at hub height (62m). The data also feature the temporal aspect, which is the time of the day, the spatial aspect, which is the coordinates, and the meteorological aspect, which includes temperature, barometric pressure, relative humidity, and wind direction. These measurements were recorded every 10 minutes for all variables. For ensuring comparability and removing scale bias: Imputing missing data - Feature-wise means were used to replace missing values; Standardisation - Numerical features are standardised to have a zero average and unit variance; Normalisation - Normalisation of spatial coordinates to remove bias in scales of latitude and longitude. This preprocessing helps stabilise the training process for both the Linear and the Gaussian Process models.

Clustering

Wind conditions in several neighbouring wind farms are likely similar to each other because they are subject to similar weather conditions. To reduce computation, several wind farms can be combined into clusters, enabling forecasts based on clusters rather than individual wind farms. There is an advantage to combining information from several wind farms in the cluster, as it may help iron out local weather extremes.

To capture regime-dependent wind dynamics, stations were clustered using geographical coordinates (longitude, latitude) and elevation. We first used the Hopkins test to assess the dataset's tendency to form clusters and determine whether the data are meaningfully clustered or randomly distributed (i.e., not clusterable). According to the method proposed by Hopkins and Skellam [18], the clustering tendency can be quantified using the Hopkins statistic H defined as:

$$H = \frac{\sum_{i=1}^n u_i^d}{\sum_{i=1}^n u_i^d + \sum_{i=1}^n w_i^d} \quad (1)$$

where $H > 0.75$ suggests strong evidence of the cluster structure, while $H < 0.5$ suggests little or no cluster tendency. Thereafter, we will use hierarchical clustering to group the data. This is a method of cluster analysis that builds a hierarchy of clusters to uncover the underlying structure in a dataset. It does not rely on probability distributions, making it a nonparametric method, suitable even when the underlying data distribution is unknown. It also produces a dendrogram, a tree diagram, that provides a visual, flexible representation of how clusters form at different levels. This allows users to choose the number of clusters by cutting the tree at the desired level.

Variable Selection - Relaxed Lasso

The Relaxed Lasso is a two-stage modification of the standard Lasso regression that aims to reduce the bias introduced by Lasso's shrinkage, particularly in high-dimensional settings [19]. Lasso performs both variable selection and shrinkage by minimising Equation 2: Lasso estimates regression coefficients by solving the Equation 2. As shown in Equation 2, the Lasso introduces an ℓ_1 -penalty to enforce sparsity.

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \left\{ \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\} \quad (2)$$

The factor $\frac{1}{2n}$ expresses the loss as average empirical risk; alternative scalings (e.g., $\frac{1}{2}$) are equivalent up to a rescaling of the tuning parameter λ . While effective for variable selection, the Lasso tends to over-shrink coefficient estimates, leading to biased results. The Relaxed Lasso addresses this by decoupling the variable selection and coefficient estimation steps. Specifically, it first uses the Lasso to select a subset of predictors and then performs a second estimation step with reduced shrinkage on those selected variables. The Relaxed Lasso estimator is defined as Equation 3:

$$\hat{\beta}^{\text{relaxed}}(\lambda, \phi) = \phi \cdot \hat{\beta}^{\text{lasso}}(\lambda) + (1 - \phi) \cdot \hat{\beta}^{\text{LS}}(\lambda), \quad (3)$$

where $\hat{\beta}^{\text{lasso}}(\lambda)$ is the Lasso solution at regularisation level λ , $\hat{\beta}^{\text{LS}}(\lambda)$ is the ordinary least squares estimate on the set of variables selected by the Lasso, $\phi \in [0, 1]$ is the relaxation parameter, controlling the degree of shrinkage. When $\phi = 1$, the relaxed estimator reduces to the standard Lasso; when $\phi = 0$, it corresponds to unregularised least squares on the Lasso-selected model. Intermediate values of ϕ provide a bias-variance trade-off.

2.2. GPR Models for Wind Forecasting

Gaussian processes are a powerful nonparametric Bayesian approach to regression and classification introduced by [20]. These processes are conceptualised as components of spatial-temporal modelling, specifically within a defined region where stations are located at different sites. Spatial analysis seeks to construct an optimal model for generating outputs by considering inputs from diverse locations over various time frames. Equation (4) describes a spatio-temporal GP model:

1. Model equation:

$$Y(s_i, t) = \mathbf{x}^\top(s_i, t)\boldsymbol{\beta} + w(s_i, t) + \varepsilon(s_i, t), \quad (4)$$

where $i = 1, 2, \dots, n$, $t = 1, 2, \dots, T$, $Y(s_i, t)$ is the observed data at location s_i and time t . and

$$w(s_i, t) \sim GP(0, C((s_i, t), (s_j, t'))), \quad (5)$$

$$\varepsilon(s_i, t) \sim \mathcal{N}(0, \sigma^2). \quad (6)$$

2. Prior to the latent process:

The independent Gaussian process (GP) model is specified hierarchically by:

$$Y_{it} = O_{it} + \varepsilon_{it}, \quad (7)$$

$$O_{it} = X_{it}\boldsymbol{\beta} + w_{it}, \quad (8)$$

for each $i = 1, \dots, n$ and $t = 1, \dots, T$, where we assume that ε_{it} and w_{it} are independent, and each is normally distributed with its respective parameters. The notation $\varepsilon_{it} = (\varepsilon(s_1, t), \dots, \varepsilon(s_n, t))^\top$ will be used to denote the so-called *nugget effect* or the pure error term, which is assumed to be independently normally distributed as

$$\varepsilon_{it} \sim \mathcal{N}(0, \sigma_\varepsilon^2 I),$$

where σ_ε^2 is the unknown pure error variance and I is the identity matrix.

The spatio-temporal random effects will be denoted by $w_{it} = (w(s_1, t), \dots, w(s_n, t))^\top$, and these will be assumed to follow

$$w_{it} \sim \mathcal{N}(0, \Sigma)$$

independently in time

Let \mathcal{O} denote all the random effects O_{it} , for $i = 1, \dots, n$ and $t = 1, \dots, T$. Let

$$\Theta = (\beta, \sigma_\varepsilon^2, \sigma_w^2, \rho, \nu)$$

denote all the parameters of this model, and let $\pi(\Theta)$ denote the prior distribution.

The logarithm of the joint posterior distribution of the parameters and the missing data for this GP model is given by:

$$\begin{aligned} \log \pi(\Theta, \mathcal{O}, y^* | y) \propto & -\frac{N}{2} \log \sigma_\varepsilon^2 - \frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^n \sum_{t=1}^T (Y_{it} - O_{it})^\top (Y_{it} - O_{it}) \\ & - \sum_{i=1}^n \frac{T}{2} \log |\sigma_w^2| - \frac{1}{2\sigma_w^2} \sum_{i=1}^n \sum_{t=1}^T (O_{it} - X_{it}\beta)^\top S_w^{-1} (O_{it} - X_{it}\beta) \\ & + \log \pi(\Theta). \end{aligned} \quad (9)$$

3. Linear models:

The very first attempt at modelling a spatio-temporal response variable is to consider the linear regression models. The multiple linear regression model is written in Equation (10) as:

$$Y(s_i, t) = \beta_1 x_1(s_i, t) + \dots + \beta_p x_p(s_i, t) + \varepsilon(s_i, t), \quad \text{for } i = 1, \dots, n; t = 1, \dots, T, \quad (10)$$

where $\varepsilon(s_i, t)$ is the spatio-temporal error term. The error term $\varepsilon(s_i, t)$ is assumed to be a zero-mean spatio-temporal Gaussian process with a separable covariance structure, given by:

$$\text{Cov}(\varepsilon(s_i, t_k), \varepsilon(s_j, t_l)) = \sigma^2 \rho_s(\|s_i - s_j\|; \theta_s) \rho_t(|t_k - t_l|; \theta_t),$$

where $\rho_s(\|s_i - s_j\|; \theta_s)$ is an anisotropic correlation function of the spatial distance $\|s_i - s_j\|$ between two locations s_i and s_j , which may depend on the parameter θ_s , which may be more than one in number.

2.3. Hybrid Model

The proposed Hybrid GP Regression combines a linear regression model with a Gaussian Process (GP) applied to residuals, enabling both interpretable trend estimation and flexible modelling of spatio-temporal dependencies.

Hybrid Model Formulation

Let $Y(\mathbf{s}_i, t)$ denote the wind speed observed at spatial location \mathbf{s} and time t . The proposed model is defined as:

$$Y(\mathbf{s}_i, t) = \mathbf{x}(\mathbf{s}_i, t)^\top \boldsymbol{\beta} + f(\mathbf{s}_i, t) + \varepsilon(s_i, t), \quad (11)$$

with

$$f(s_i, t) \sim GP(0, k((s_i, t), (s_j, t'))), \quad (12)$$

and

$$\varepsilon(s_i, t) \sim \mathcal{N}(0, \sigma^2). \quad (13)$$

As shown in Equation (11), the model combines a linear regression component with a Gaussian process for residuals. This two-stage formulation enables a clear separation between deterministic structure and stochastic variability. The linear component captures large-scale spatio-temporal patterns using meteorological and temporal covariates. Model parameters are estimated using ordinary least squares. The fitted linear model serves as a baseline and provides residuals for subsequent GP modelling.

Gaussian Process Residual Model

The residual process ($f(\mathbf{s}, t)$) is modelled using a zero-mean GP with covariance function (k). Multiple covariance structures are considered, including: Matérn covariance with smoothness parameters ($\nu = 0.5, 1.5, 2.5$), Exponential covariance, and Gaussian (squared exponential) covariance. Each covariance structure encodes different assumptions about the smoothness of the residual wind-speed field.

Covariance Selection Procedure

For every validation split of residuals, GP models with different covariance functions are employed. Estimation of the hyperparameters is conducted through maximum likelihood optimisation. When selecting the best covariance function, a choice is made based on validation RMSE and probabilistic performance measures.

Validation Strategy and Evaluation Metrics

A leave-one-station-out cross-validation method was used to evaluate the model's performance. In each iteration, validation was performed on one station, and training was conducted on the remaining stations. This repetition continued until all stations were used for validation once. Hence, for four stations, 25% of the data was utilised for validation; for six stations, 16.7% was utilised. This checks the model's ability to make predictions at unobserved locations of wind measurement.

2.4. Validation Metrics

Bayesian evaluation metrics are used to evaluate the performance of Bayesian models, such as Bayesian neural networks. These metrics often go beyond traditional accuracy measures to incorporate uncertainty estimates, which are a key feature of Bayesian methods. The following techniques are used to evaluate the models and are classified as deterministic and probabilistic. The deterministic metrics are the Root Mean Square Error (RMSE) and Mean Absolute Error (MAE), and the probabilistic metrics are the continuous ranked probability score (CRPS), the coverage probability (CVG), the negative log predictive density (NLPD), PINAD and PINAW.

RMSE (Root Mean Squared Error):

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}, \quad (14)$$

MAE (Mean absolute error):

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|, \quad (15)$$

where y_i is the wind speed, \hat{y}_i is the forecast of the wind speed, and N is the number of forecasts. RMSE measures the average magnitude of forecast errors, giving a higher weight to larger errors, whereas MAE measures the average absolute difference between predicted and actual values.

The probabilistic metrics are given below:

$$\text{CRPS}(F, y) = \int_{-\infty}^{\infty} (F(z) - \mathbf{1}\{y \leq z\})^2 dz \quad (16)$$

The Continuous Ranked Probability Score (CRPS) can be written in expectation form [21] as

$$\text{CRPS}(F, y) = \mathbb{E}_F[|Y - y|] - \frac{1}{2} \mathbb{E}_F[|Y - Y'|], \quad Y, Y' \sim F \text{ i.i.d.} \quad (17)$$

where: F is the predictive cumulative distribution function (CDF), y is the observed value $Y \sim F$, $Y' \sim F$ are independent random variables drawn from the forecast distribution F and \mathbb{E}_F denotes

the expectation under the distribution F . CRPS evaluates the quality of a probabilistic forecast by comparing the forecast's cumulative distribution function (CDF) with the observed value. CVG measures the proportion of true observations that fall within a predicted confidence interval of a given nominal level (e.g., 95%). The predictive coverage of the $1 - \alpha$ confidence interval is measured using the Coverage (CVG) metric. For n observations, let \hat{y}_i^L and \hat{y}_i^U denote the lower and upper bounds of the predicted $(1 - \alpha) \times 100\%$ confidence interval for the i th observation. The empirical coverage is then

$$\text{CVG}_{1-\alpha} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i \in [\hat{y}_i^L, \hat{y}_i^U]\}, \quad (18)$$

where y_i is the observed value at test point i , $\mathbf{1}[\cdot]$ is the indicator function, returning 1 if the condition is true and 0 otherwise. CVG represents the proportion of observations falling within their corresponding predictive intervals. A well-calibrated model should achieve $\text{CVG}_{1-\alpha} \approx 1 - \alpha$, reflecting accurate uncertainty quantification.

Negative Log Predictive Density (NLPD) to evaluate uncertainty calibration. For a single forecast-observation pair:

$$\text{NLPD} = -\log p(y | \mathcal{D}), \quad (19)$$

where y is the observed value, $p(y | \mathcal{D})$ is the predictive probability density at y given training data \mathcal{D} .

For N test points, the average NLPD is:

$$\text{NLPD}_{\text{avg}} = -\frac{1}{N} \sum_{i=1}^N \log p(y_i | \mathcal{D}) \quad (20)$$

The negative log predictive density (NLPD) measures how likely the model thinks the observed data are, accounting for both the mean and the uncertainty. The Prediction Interval Normalised Average Width (PINAW) measures the average width of the prediction interval, normalised by the range of observed values. Narrower intervals indicate higher precision. It is defined as:

$$\text{PINAW} = \frac{1}{k \cdot R} \sum_{i=1}^k (U_i - L_i), \quad (21)$$

where: k is the number of validation points, U_i and L_i are the upper and lower bounds of the prediction interval for the i -th observation, $R = \max(y_{\text{true}}) - \min(y_{\text{true}})$ is the range of the observed values.

The Prediction Interval Normalised Average Deviation (PINAD) measures the average deviation of the midpoint of the prediction interval from the true values, normalised by the range. It indicates the calibration of the prediction intervals. It is defined as:

$$\text{PINAD} = \frac{1}{k \cdot R} \sum_{i=1}^k \left| \frac{U_i + L_i}{2} - y_{\text{true},i} \right|, \quad (22)$$

where the symbols have the same meaning as above, lower values of PINAD indicate that the prediction interval is well centred around the observed values.

3. Results

In this section, the exploratory data analysis will be outlined and evaluated. R statistical software version 4.5.2 was used for data analysis.

3.1. Exploratory Data Analysis (EDA)

As part of analysing wind speed data, Exploratory Data Analysis (EDA) uses statistical measures and various forms of data visualisation. EDA, in particular, becomes important when forming spatio-

temporal data models, as it reveals inherent data trends, structures, and patterns, providing the foundation for selecting an ideal model and specifying parameters.

3.1.1. Clustering

In spatio-temporal predictions, clustering the locations balances efficiency, accuracy, and interpretability. This makes clustering a powerful strategy for large-scale forecasting tasks such as weather, traffic, and energy production.

Hopkins Test

The Hopkins Statistic is a method for evaluating whether a dataset has a meaningful clustering structure, i.e., its tendency to form clusters. It is often used before applying clustering algorithms to spatial data. The Hopkins test is used to assess the tendency for clustering in a dataset.

- **Null Hypothesis (H_0):** The data is uniformly randomly distributed in the feature space (i.e., there is no meaningful cluster structure).
- **Alternative Hypothesis (H_1):** The data is not uniformly random and exhibits a clustering tendency.

Decision Rule:

- If the Hopkins statistic $H \approx 0.5$, we fail to reject H_0 — the data is likely random.
- If $H \gg 0.5$ (typically $H > 0.75$), we reject H_0 — the data shows a significant cluster tendency.

Since the Hopkins statistic is $H = 0.6680092 \gg 0.5$, we reject the null hypothesis and conclude that the data is clusterable.

Clustering

Since we established that our data is clusterable, the next step is to estimate the optimal or maximum number of clusters that can be meaningfully formed. The silhouette measures how well each point fits within its own cluster relative to other points. It provides a more quantitative evaluation and suggests that 2 clusters can be used for the data. Table 2 shows the stations grouped into two clusters according to their geographic location and elevation.

Table 2. Station data with clusters.

Stn no	StationCode	StationName	lon	lat	elev	cluster
1	WM1	Alexander Bay	16.66441	28.60188	152	1
2	WM2	Calvinia	19.36075	31.52494	824	1
3	WM3	Vredendal	18.41992	31.73051	242	1
4	WM5	Napier	19.69245	34.61192	288	1
5	WM6	Sutherland	20.69124	32.55680	1581	2
6	WM7	Beaufort West	22.55667	32.96672	1047	2
7	WM8	Humansdorp	24.51436	34.10997	110	1
8	WM9	Noupoort	25.02838	31.25254	1806	2
9	WM12	Eston	30.52871	29.85026	770	2
10	WM19	Upington	20.56833	27.72670	848	1

Spatial Distribution of Clusters

The scatter plot Figure 1 illustrates the spatial distribution of clusters based on station coordinates in terms of latitude and longitude. The two clusters are visually distinct—Cluster 1 (shown in red) is primarily located in the western region with lower longitude values. In contrast, Cluster 2 (shown in blue) is positioned toward the eastern part of the study area. Lastly, the coverage of the spread of Cluster 1 stations reveals a rather wide latitude range, demonstrating greater dispersion, unlike the Cluster 2 stations, which are few and geographically concentrated. This spatial differentiation is reflected in the clusters obtained, underscoring the importance of the geographical factor in their formation and highlighting climatic differences among the WASA stations.

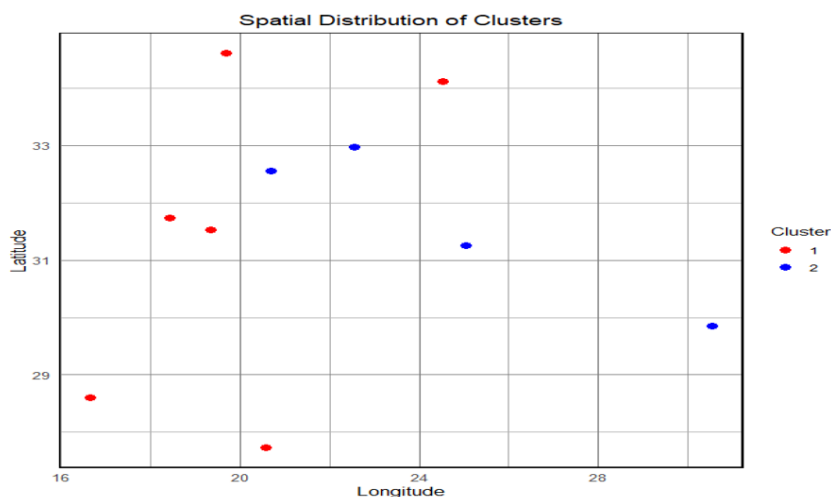


Figure 1. Spatial distribution of station clusters.

Figure 2 displays the spatial and topographical characteristics of two clusters based on their longitude, latitude, and elevation status. The longitude plot shows that stations in Cluster 1 are primarily in a lower-latitude zone, whereas those in Cluster 2 are located farther east, with higher longitudes. The plot showing the distribution based on latitude shows that while Cluster 1 covers a broader zone, extending further into the northern latitudes, Cluster 2 stations are concentrated within a narrower zone of lower latitudes. Finally, based on elevation, a notable difference between the two clusters is observed in Figure 2. In this regard, stations within Cluster 1 are concentrated in lower-elevation zones, whereas those within Cluster 2 are concentrated in higher-elevation zones. Overall, Figure 2 shows notable differences between clusters, thus supporting a spatial coherence between clusters based on longitude, latitude, and elevation status.

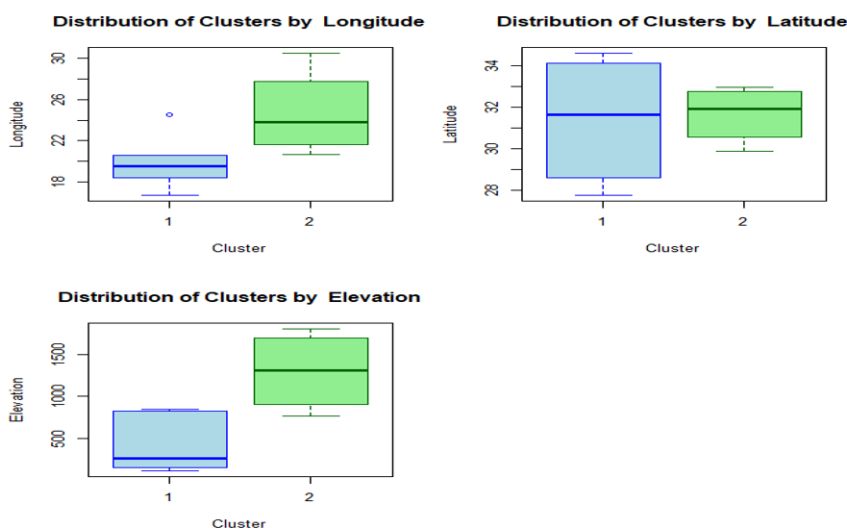


Figure 2. Distribution of clusters by coordinates and elevation.

Figure 3 includes two clusters, 1 and 2, with each cluster included in a polygon that denotes the cluster boundary. Here, Dim1 accounts for 43.7% of the variance, while Dim2 accounts for 34%. These values add up to over 77%. Thus, it is evident that this plot includes most of the clusters. In Figure 3, it is evident that the spatial distribution of the clusters is demonstrated through hierarchical methods. All the clusters generated by the dendrogram are clearly separated, demonstrating that hierarchical

clustering effectively distinguishes between them. In this plot, it is evident that the clusters are well separated in Dim1, indicating that it is an important dimension for distinguishing between them.

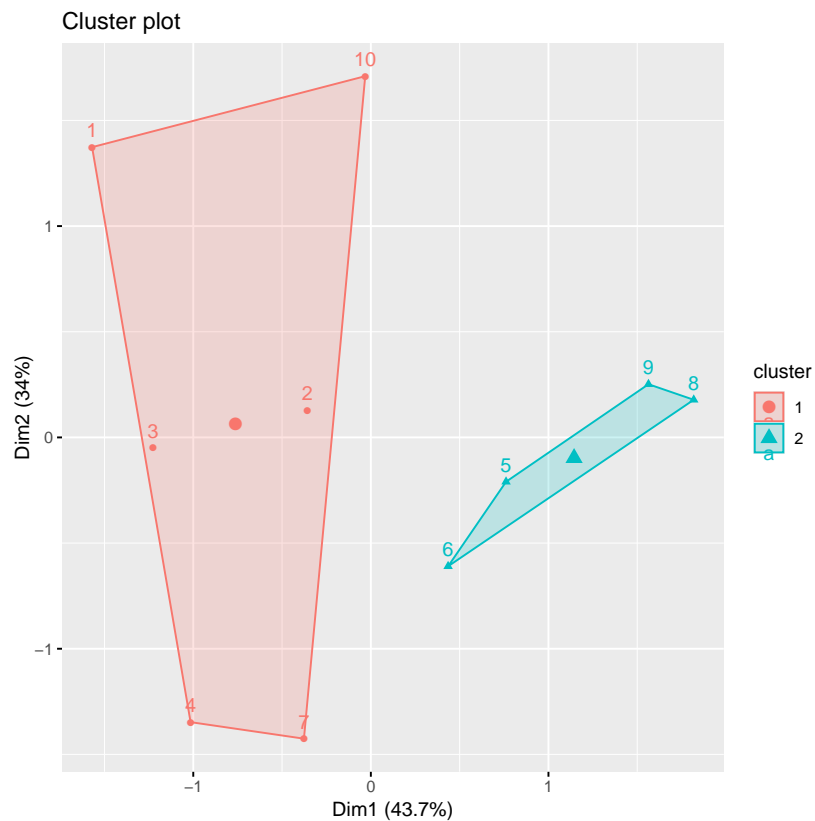


Figure 3. Hierarchical cluster plot.

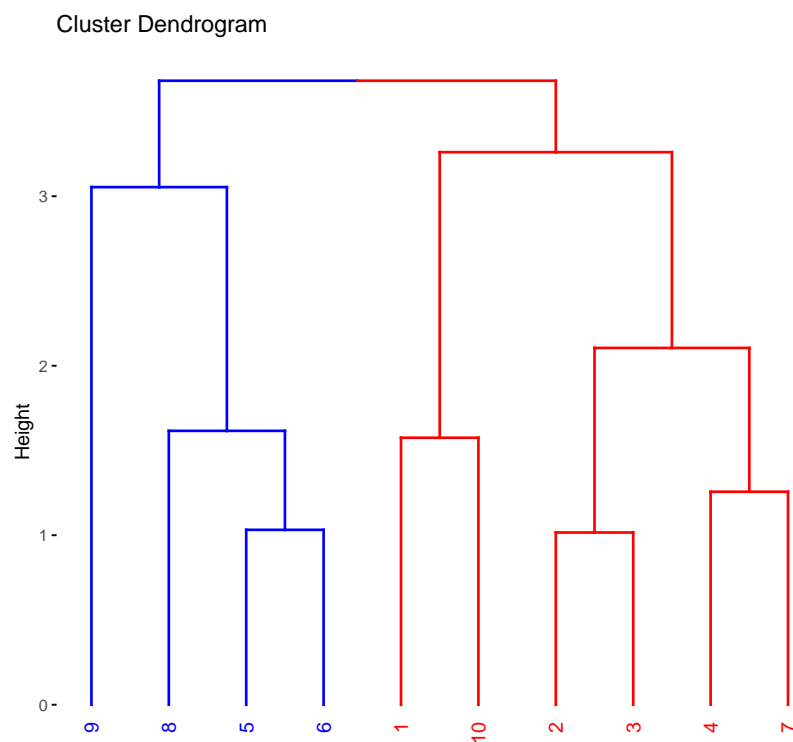


Figure 4. The dendrogram.

Figure 4 shows a dendrogram for clustering ten different stations labelled from 1 to 10. The vertical axis is labelled “Height,” indicating the degree of dissimilarity or the distance at which clustering occurs. The degree of dissimilarity is higher for points clustered at higher positions than for those clustered at lower positions. From the given dendrogram, two clusters are marked within dashed boxes. The first cluster consists of stations 9, 8, 5, and 6, whereas the second cluster consists of stations 1, 10, 2, 3, 4, and 7.

Multivariate Exploratory Data Analysis

Table 3 lists the summary statistics of wind data across stations. Many stations exhibit moderate to high variability, with median wind speeds typically 6–7.5 m/s. The highest wind activity comes from Stations 3 and 4, while Station 9 is the calmest. The maximum wind speed is close to 24 m/s, suggesting occasional high-wind events.

Also, the average values across the stations are quite different as indicated in Table 3. Alternatively, for Station 1 (stn1), the mean is higher than the median, suggesting the data may be right-skewed. Also, the difference between Q3 and the maximum is very high, suggesting the presence of outliers or a long tail.

Table 3. Summary statistics for sensor readings (stn1 to stn10).

Statistic	stn1	stn2	stn3	stn4	stn5	stn6	stn7	stn8	stn9	stn10
Min	0.3818	0.2533	0.3869	0.2075	0.2406	0.6305	0.2344	0.4235	0.2112	0.6037
1st Qu.	4.4349	3.3906	4.7480	4.9616	4.4436	4.6944	4.8155	5.3201	3.1822	4.4880
Median	7.1031	5.4118	7.5978	7.5278	6.1062	7.1340	6.7533	7.4886	4.8019	6.0382
Mean	7.5087	5.8444	7.8269	7.9377	6.6211	6.9727	7.0867	7.8101	5.0163	6.1670
3rd Qu.	9.9906	7.5868	10.7213	10.9338	8.4525	8.9410	9.2498	10.0568	6.4848	7.6730
Max	21.6708	20.5868	23.9921	20.3895	18.5664	17.7714	20.0661	19.3474	18.4819	16.5804

Figure 5 is a box plot showing the distribution of wind speed (in m/s) across 10 stations (stn1 to stn10). Each box shows the central 50% of the data (interquartile range), the line within the box shows the median wind speed, and the dots represent outliers. Stations such as stn4, stn6, and stn7 exhibit higher wind speeds and greater variability, while others, such as stn1, stn2, and stn8, have lower, more consistent wind speeds.

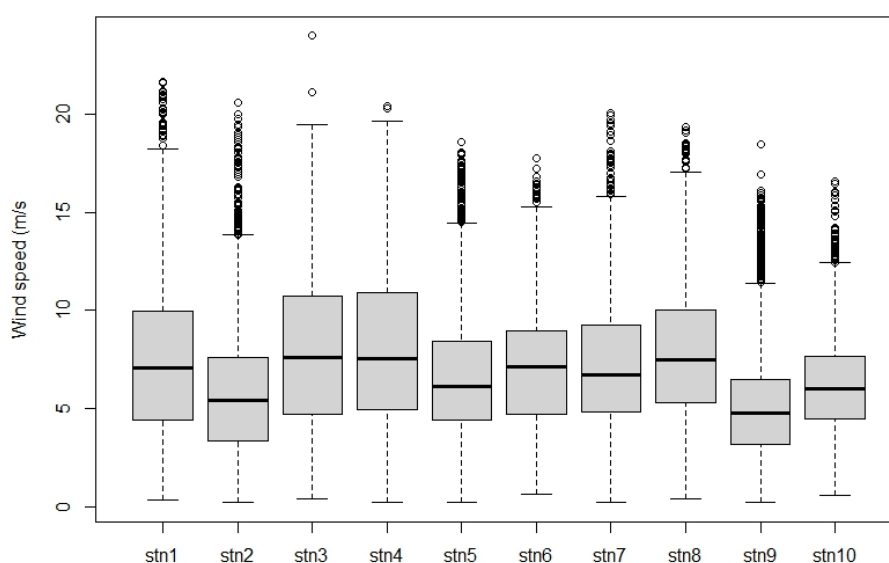


Figure 5. Boxplot of the distribution of wind speed (in m/s) across all stations.

The plot figure 6 is a time series line plot showing the wind speed (in m/s) recorded at 10 different stations over a series of observations. A sample of 500 was taken to improve visualisation of wind speed. Each line represents the variation in wind speed at a specific station, showing temporal fluctuations and variability. The overlapping lines illustrate how wind speed trends differ and coincide between stations during the observed period.

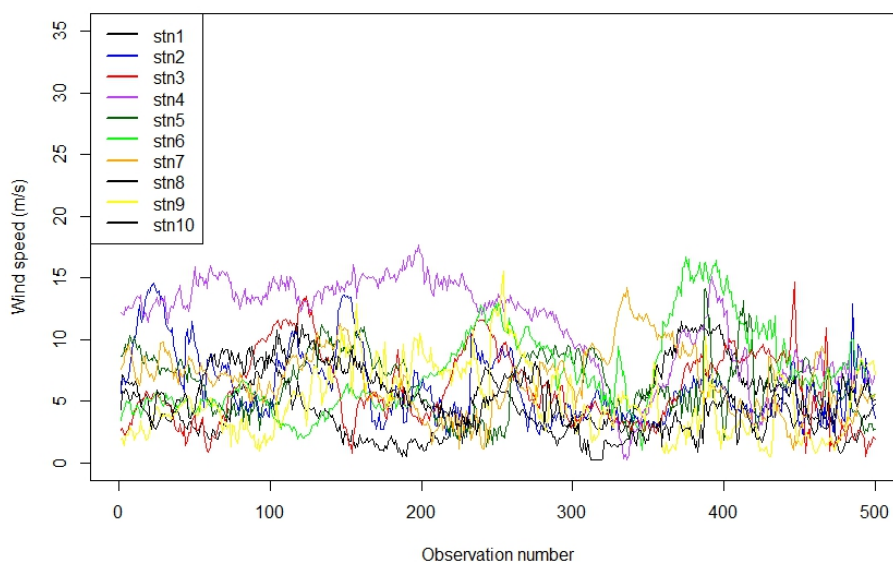


Figure 6. Time series plot for all the stations.

Figure 7 shows a scatterplot matrix combined with histograms and correlation coefficients for wind speed data. This plot is useful for understanding the pairwise relationships and dependencies among wind speed measurements across different stations. The histograms show the wind speed distribution for each station. The scatterplots show pairwise relationships between wind speeds at different stations. Each plot has a fitted smooth curve (blue line) to show trends. Correlation coefficients between wind speeds at pairs of stations are also presented. These values quantify the strength and direction of linear relationships (positive, negative, or near zero). The correlation coefficients are mostly positive but low to moderate (generally below 0.4), indicating a weak linear association. Stations such as stn1 and stn2, or stn2 and stn3, have somewhat higher correlations.

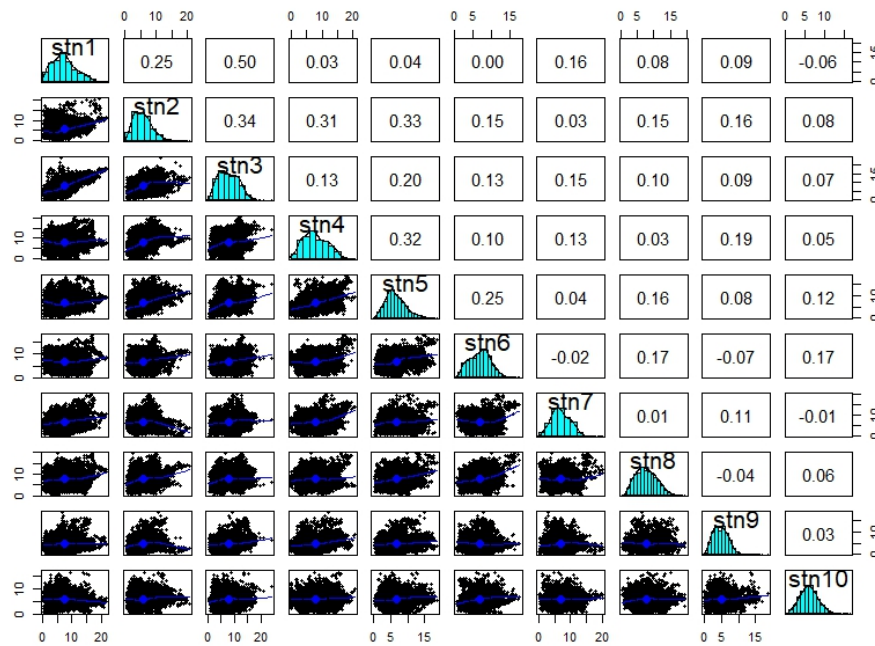


Figure 7. Correlation plot for all the stations.

Figure 8 gives a clear visualisation of how wind speed is distributed at each station. Most stations experience wind speeds typically between 0 and 15 m/s, with varying frequencies. Differences in histogram shape suggest variations in local wind-speed patterns across stations.

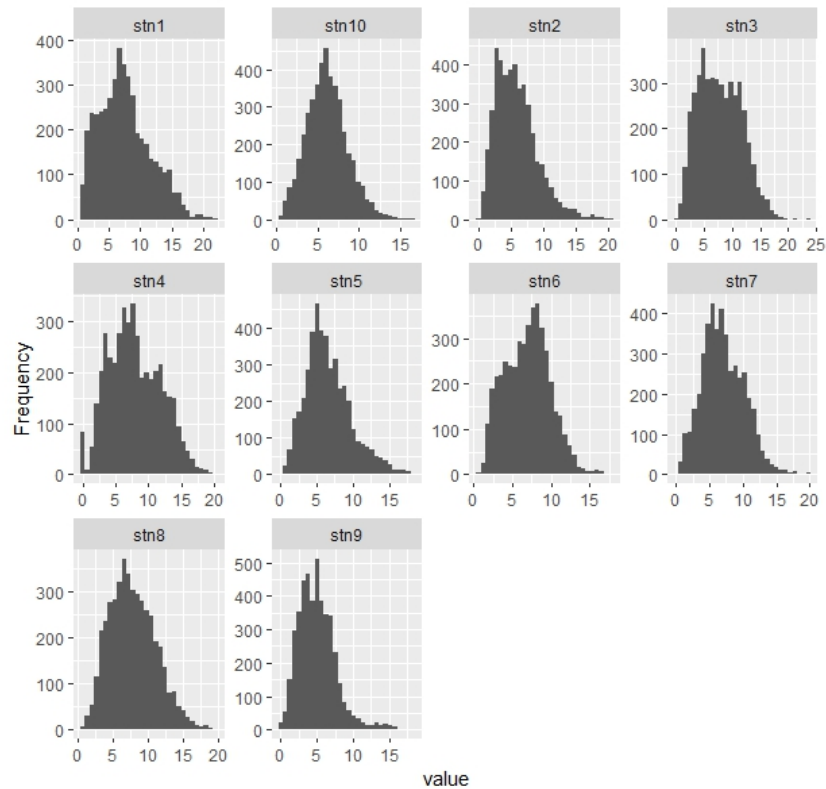


Figure 8. Histogram of all the stations.

Figure 9 shows the correlation between stations. From the plot, stations 1 and 2 show a strong positive relationship. As the wind speed of station 1 increases, so does that of station 3.

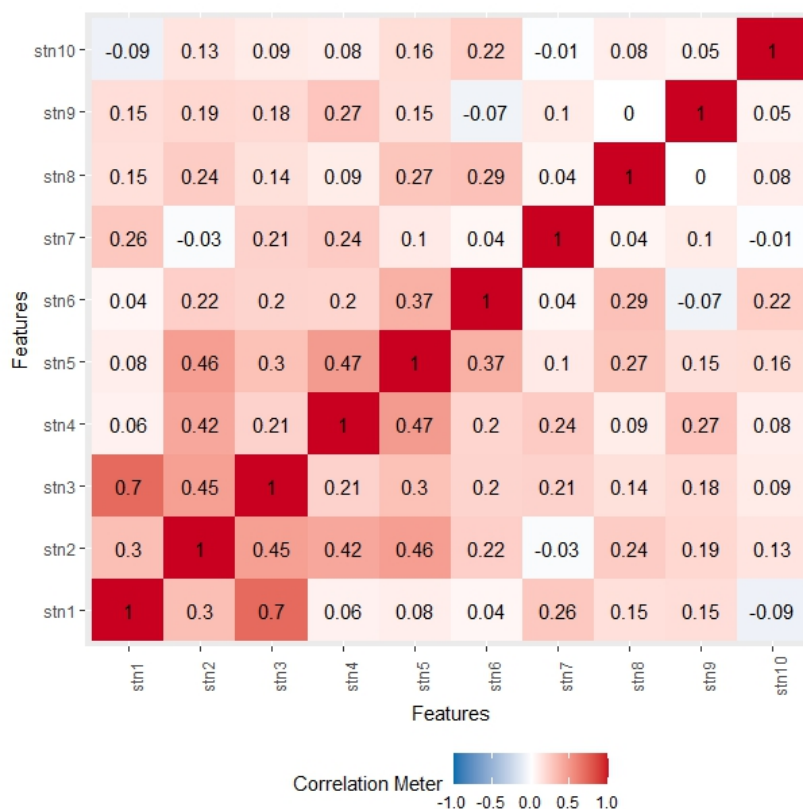


Figure 9. Correlation plot.

3.2. Results for the Models

The results for the models are presented in the following sections in two stages. Firstly, section 3.2.1 presents the results for Gaussian Process (GP) models applied to all possible combinations of train-test stations, providing a comprehensive evaluation of predictive performance and metrics such as RMSE, NLPD, CVG, and acceptance rate ϕ . These results establish the baseline behavior of the GP models across clusters and highlight the effects of station selection and spatial coverage.

In contrast, Section 3.2.2 evaluates the hybrid Linear–Gaussian Process framework. Here, a linear regression model is first estimated to capture large-scale spatio-temporal trends, and a GP is subsequently applied to the residuals to model local deviations. Within this second stage, a variety of covariance structures are rigorously assessed, including Matérn functions with smoothness parameters $\nu = 0.5, 1.5, 2.5$, as well as exponential and Gaussian covariances. The final covariance choice is determined based on validation performance. Because the hybrid model explicitly separates linear and residual components, the train–test splits, validation folds, and reported metrics differ from Section 3.2.1. Metrics such as ϕ are not always relevant for residual modeling, and certain summary metrics (e.g., CRPS) are computed only where meaningful for the hybrid framework. This structure ensures that the evaluation in Section 3.2.2 reflects the performance of the hybrid approach in capturing both large-scale trends and local residuals, rather than replicating the exhaustive GP-only analysis.

3.2.1. Results for the Linear and GP Models

Variable Selection

Relaxed Lasso is a variation of the Lasso (Least Absolute Shrinkage and Selection Operator) regression method, designed to improve variable selection by reducing the bias introduced by Lasso's strong penalty. In the first step, Lasso is applied to select a subset of relevant variables by shrinking some coefficients exactly to zero. Then, in the second step, a less-penalised or unpenalised regression (like OLS) is fit only on the selected variables, allowing the model to estimate their coefficients more

accurately. This two-stage process helps maintain Lasso's variable selection strength while improving the accuracy and interpretability of the prediction by relaxing the amount of shrinkage applied.

Forecast Evaluation

We perform GP modelling using the spTimer package in R, and we use an MCMC approach to estimate the posterior distributions of the model parameters. The following MCMC settings were used: number of iterations (N): 2,000; burn-in period: 300 iterations; thinning interval: 1; number of chains: 3; and reporting interval (report): 3. Such settings mean that the program was run for 2,000 iterations, of which the initial 300 were considered as burn-in to minimise the impact of initial values. Thinning of 1 means that the entire set of post-burn samples is utilised for posterior analysis. Running the chains enables convergence diagnosis, ensuring proper parameter evaluation. The report = 3 settings ensure that the algorithm reports every 3rd iteration, thus enabling tracking of the execution of the sampling function. Such settings ensure an appropriate trade-off between computational efficiency and the sample size required to approximate the posterior distribution.

Evaluation for cluster 1

Training and testing for cluster 1

In cluster 1, we use the ratio $\frac{6}{2} = 3$ to generate all possible combinations of data for the test set, with test set sizes ranging from 1 to 3 samples. We used the configurations $\binom{6}{1} = 6$, $\binom{6}{2} = 15$ and $\binom{6}{3} = 20$, which resulted in 41 unique train-test partitions. In each configuration, the selected subset of stations served as the test set and the remaining stations were used for training. The objective of this exhaustive enumeration is to evaluate the robustness of spatial generalisation with respect to the specific choice of held-out stations. By considering every possible hold-out combination, we eliminate randomness due to fold assignment and provide a complete assessment of how sensitive model performance is to station selection. The combinations resulting from the testing and training sets are in Table 4.

Table 4. Testing and Training Set Combinations (Grouped into 6C1–6C2 and 6C3).

6C1 and 6C2			6C3		
ID	Testing Set	Training Set	ID	Testing Set	Training Set
1	1	2, 3, 4, 5, 6	22	1, 2, 3	4, 5, 6
2	2	1, 3, 4, 5, 6	23	1, 2, 4	3, 5, 6
3	3	1, 2, 4, 5, 6	24	1, 2, 5	3, 4, 6
4	4	1, 2, 3, 5, 6	25	1, 2, 6	3, 4, 5
5	5	1, 2, 3, 4, 6	26	1, 3, 4	2, 5, 6
6	6	1, 2, 3, 4, 5	27	1, 3, 5	2, 4, 6
7	1, 2	3, 4, 5, 6	28	1, 3, 6	2, 4, 5
8	1, 3	2, 4, 5, 6	29	1, 4, 5	2, 3, 6
9	1, 4	2, 3, 5, 6	30	1, 4, 6	2, 3, 5
10	1, 5	2, 3, 4, 6	31	1, 5, 6	2, 3, 4
11	1, 6	2, 3, 4, 5	32	2, 3, 4	1, 5, 6
12	2, 3	1, 4, 5, 6	33	2, 3, 5	1, 4, 6
13	2, 4	1, 3, 5, 6	34	2, 3, 6	1, 4, 5
14	2, 5	1, 3, 4, 6	35	2, 4, 5	1, 3, 6
15	2, 6	1, 3, 4, 5	36	2, 4, 6	1, 3, 5
16	3, 4	1, 2, 5, 6	37	2, 5, 6	1, 3, 4
17	3, 5	1, 2, 4, 6	38	3, 4, 5	1, 2, 6
18	3, 6	1, 2, 4, 5	39	3, 4, 6	1, 2, 5
19	4, 5	1, 2, 3, 6	40	3, 5, 6	1, 2, 4
20	4, 6	1, 2, 3, 5	41	4, 5, 6	1, 2, 3
21	5, 6	1, 2, 3, 4			

Results for Cluster 1

Table 5 summarizes predictive performance across all 41 spatial hold-out configurations. When a single station is held out (leave-one-station-out validation), the model is trained on five stations and evaluated on one. As expected, this setting yields the lowest RMSE, MAE, NLPD, and CRPS values because the training set is largest and the extrapolation domain is minimal. When two stations are held out, predictive performance degrades moderately across all proper scoring rules. This reflects both the reduced amount of spatial information available for training and the increased geographic extent of the test domain. The most challenging setting occurs when three stations are held out. In this case, the model is trained on only half of the available stations and must extrapolate to a substantially larger unseen spatial region. Consequently, RMSE, MAE, CRPS, and NLPD increase consistently in relation to cases of one and two-stations.

This monotonic degradation pattern indicates that performance is sensitive to the degree of spatial data removal rather than to any specific station combination. Importantly, coverage (CVG) remains stable across configurations, suggesting that predictive uncertainty is well calibrated even as extrapolation difficulty increases. The variability observed within each group (for example, across different three-station subsets) reflects heterogeneity in spatial information content across stations, highlighting that some station combinations are inherently more informative for spatial interpolation than others. Therefore, Table 5 quantifies how predictive accuracy deteriorates as the spatial extrapolation task becomes more challenging. In this sense, the table illustrates the sensitivity of model performance to the degree of spatial data removal, providing a structured robustness assessment of spatial extrapolation performance. The metric ϕ quantifies the relative contribution of each train-test combination within the fusion scheme. Specifically, for combination j , ϕ_j is defined as

$$\phi_j = 100 \times \frac{w_j}{\sum_k w_k}, \quad (23)$$

where w_j represents the weight or influence assigned to combination j based on its predictive performance. CVG and other metrics are used to determine w_j in a manner that emphasizes well-performing combinations. The resulting ϕ_j values are expressed as percentages, highlighting which combinations contribute the most to the overall fused prediction. It represents the acceptance rate in the GPR fusion scheme, highlighting the train-test combinations that contribute most reliably to the final predictive ensemble. The metric ϕ further highlights the combinations that contribute the most within the fusion scheme, with certain station sets such as 2, 3, 3-5, 3-6, 2-5-6, receiving higher relevance scores, suggesting that these combinations provide more informative predictors for the model.

Table 5. Performance metrics for different station combinations (compact view).

Combination	NLPD	PINAW	PINAD	RMSE	MAE	CRPS	CVG	ϕ (%)
1	1.9772	0.2690	-0.0156	1.4581	1.1639	0.8857	0.9500	22.00
2	1.8799	0.2134	0.0367	1.2051	0.9906	0.7701	0.9500	92.05
3	1.8407	0.1984	-0.0194	1.2138	0.9260	0.7491	0.9500	67.55
4	2.0453	0.3163	-0.0301	1.6529	1.3192	0.9794	0.9500	22.35
5	1.9173	0.2382	0.0035	1.2164	0.9596	0.7881	0.9500	26.35
6	1.9280	0.3274	0.0339	1.3307	1.0546	0.8254	0.9500	24.35
1-2	2.1802	0.2640	0.0161	1.4966	1.2269	1.0074	0.9500	41.44
1-3	2.3404	0.4258	-0.0366	2.4553	1.9306	1.3918	0.9500	23.15
1-4	2.2288	0.3677	-0.0321	1.9697	1.5561	1.1681	0.9500	23.05
1-5	2.2467	0.3432	-0.0107	1.8504	1.4508	1.1399	0.9500	20.25
1-6	2.1970	0.3022	0.0059	1.6391	1.3080	1.0537	0.9500	18.50
2-3	2.1310	0.2640	0.0088	1.6682	1.3634	1.0251	0.9500	18.45
2-4	2.3002	0.4221	0.0162	2.2910	1.8898	1.3245	0.9500	20.40
2-5	2.1372	0.2512	0.0273	1.3396	1.0761	0.9412	0.9500	32.80
2-6	2.2779	0.3443	0.0768	1.9823	1.7042	1.2087	0.9500	32.35
3-4	2.1903	0.3035	-0.0388	1.8522	1.4662	1.1113	0.9500	53.75
3-5	2.1291	0.2397	-0.0206	1.3993	1.0840	0.9480	0.9500	36.75
3-6	2.1276	0.2310	0.0016	1.4057	1.0887	0.9475	0.9500	40.00
4-5	2.4069	0.5244	-0.0276	2.6420	2.1088	1.5004	0.9500	23.25
4-6	2.3480	0.4647	-0.0063	2.4049	1.9453	1.3858	0.9500	19.20
5-6	2.3853	0.5165	0.0261	2.5984	2.0662	1.4700	0.9500	26.70
1-2-3	2.6265	0.5341	0.0010	3.2953	2.6461	1.8635	0.9500	25.95
1-2-4	2.4450	0.4465	-0.0037	2.4452	1.9569	1.4516	0.9500	20.55
1-2-5	2.4384	0.3841	0.0179	2.1756	1.7505	1.3661	0.9500	18.95
1-2-6	2.4819	0.4118	0.0482	2.3775	1.9584	1.4631	0.9500	19.45
1-3-4	2.5240	0.5029	-0.0570	2.9707	2.3361	1.6811	0.9500	23.30
1-3-5	2.5704	0.5219	-0.0346	3.0764	2.4197	1.7442	0.9500	23.60
1-3-6	2.4954	0.4413	-0.0099	2.5856	2.0518	1.5300	0.9500	22.00
1-4-5	2.4914	0.4787	-0.0374	2.5171	1.9611	1.5019	0.9500	21.15
1-4-6	2.4247	0.3955	-0.0178	2.1086	1.6635	1.3356	0.9500	20.00
1-5-6	2.4990	0.4158	0.0050	2.2558	1.7919	1.4351	0.9500	26.40
2-3-4	2.4937	0.4780	-0.0137	2.9242	2.3500	1.6535	0.9500	20.95
2-3-5	2.3653	0.3322	0.0066	2.0923	1.6908	1.2898	0.9500	18.45
2-3-6	2.4715	0.4477	0.0390	2.8127	2.3023	1.6068	0.9500	24.30
2-4-5	2.5664	0.5905	0.0107	3.0813	2.5013	1.7589	0.9500	26.65
2-4-6	2.7089	0.6601	0.0401	3.5722	2.9373	2.0517	0.9500	25.00
2-5-6	2.5753	0.5879	0.0690	3.1232	2.5483	1.7855	0.9500	27.80
3-4-5	2.4724	0.4564	-0.0501	2.6810	2.0913	1.5440	0.9500	21.00
3-4-6	2.4396	0.3782	-0.0241	2.2505	1.7764	1.3871	0.9500	19.40
3-5-6	2.4585	0.4047	-0.0005	2.4462	1.9551	1.4616	0.9500	24.65
4-5-6	2.6000	0.5781	-0.0017	2.9612	2.3805	1.7331	0.9500	27.75

Prediction Plots

As the combined station counts increase, accuracy gradually deteriorates. At the same time, pairwise combinations tend to report higher NLPD, wider prediction intervals, and greater RMSEs as model complexity increases, due to greater spatial variability. This phenomenon is reinforced when visualised within Figure 10, where combined stations report wider bands of conformal prediction, less smooth trends for the mean trajectories, while the prediction intervals report perfectly well-calibrated performance across the different combinations, reporting the desired true coverage of 95% as affirmed by the CVG values reported within Table 5.

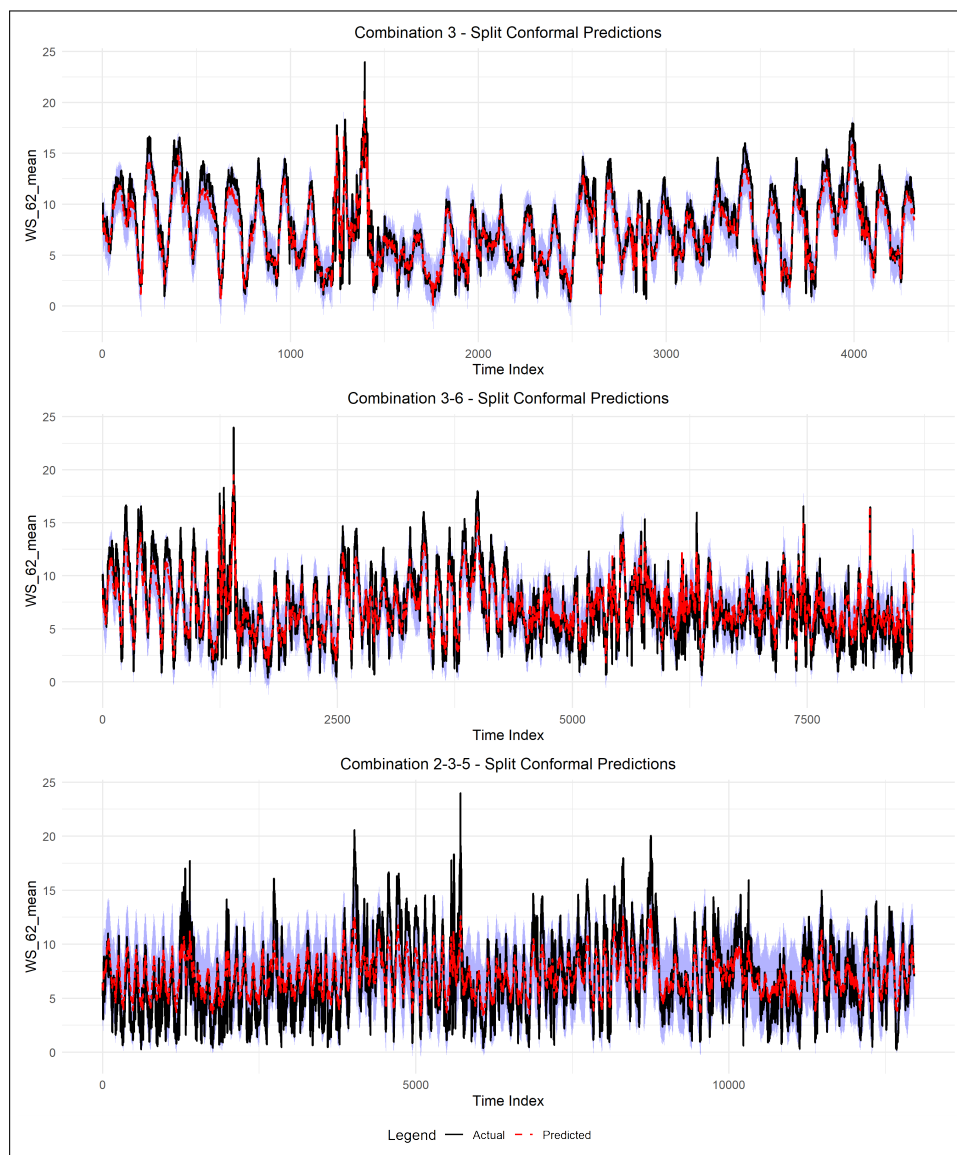


Figure 10. Prediction plots for cluster 1 high performing stations.

Importantly, the adaptive extension of the prediction interval in the regime of high wind variability, evident from Figure 10, shows the capability of the split conformal method in handling the dynamic uncertainty without sacrificing the guarantees. These examples reveal an interesting balance between the sharpness of predictions and the level of aggregation, indicating that although aggregation can be valuable, it might also entail a penalty in the loss of local wind dynamics, which is of principal importance.

Evaluation for cluster 2

Testing and Training Sets

The training set is the data used to train the model, allowing it to learn patterns and relationships. In contrast, the testing set is a separate portion of data used to assess how well the model performs on new, unseen inputs. By splitting data into these two sets, we can ensure the model learns effectively without simply memorising the data, helping to avoid overfitting and giving a more accurate measure of how the model will perform in real-world scenarios. Table 6 shows the test-train sets for cluster 2.

Table 6. Testing and Training Set Combinations for cluster 2.

No.	Testing Set	Training Set
1	1	2, 3, 4
2	2	1, 3, 4
3	3	1, 2, 4
4	4	1, 2, 3
5	1, 2	3, 4
6	1, 3	2, 4
7	1, 4	2, 3
8	2, 3	1, 4
9	2, 4	1, 3
10	3, 4	1, 2

To evaluate the model's performance and generalisation, two data-splitting strategies were employed, each using combinations of four spatial locations. We used the two configurations $\binom{4}{1} = 4$, $\binom{4}{2} = 6$, which denote combinations in which subsets of four components are split into training and testing sets. These combinations provide a systematic way to assess the model's robustness under varying data availability scenarios and to evaluate predictive performance across different spatial partitions.

Results for Cluster 2

Table 7 presents a comparative evaluation of the performance of the model in various training and testing configurations using key statistical metrics: Root Mean Square Error (RMSE), Mean Absolute Error (MAE), Continuous Ranked Probability Score (CRPS), probability of prediction interval coverage (PICP), Negative Log Predictive Density (NLPD), normalised average width of prediction interval (PINAW) and normalised average deviation of prediction interval (PINAD).

Table 7. Model performance for LM and GPR across station combinations.

s.index	Model	RMSE	MAE	CRPS	PICP	NLPD	PINAW
1	LM	0.8198	0.6321	0.4530	0.9454	1.2241	0.1728
	GP	1.1313	0.8616	1.3826	0.999	1.8957	0.507
2	LM	0.8283	0.6421	0.4585	0.9474	1.2346	0.1852
	GP	1.2536	1.003	1.3806	0.999	1.9227	0.5429
3	LM	0.8199	0.6212	0.4495	0.9379	1.2241	0.1658
	GP	2.1206	1.6249	1.2953	0.9506	2.1743	0.4679
4	LM	1.2005	0.7890	0.5963	0.8946	1.8832	0.1708
	GP	2.8746	2.4061	1.3586	0.9215	2.4990	0.5225
1,2	LM	0.8202	0.6300	0.4521	0.9463	1.2245	0.1732
	GP	2.4598	1.9836	1.5065	0.9024	2.6392	0.5685
1,3	LM	1.5695	1.3962	1.0527	0.5676	2.8405	0.1596
	GP	2.1250	2.6561	1.9941	0.8330	2.9306	0.5220
1,4	LM	1.1129	0.7302	0.5433	0.9166	1.7300	1.1689
	GP	2.7990	2.2896	1.7487	0.9099	2.7441	0.6694
2,3	LM	0.8231	0.6303	0.4536	0.9420	1.2285	1.6688
	GP	3.1901	2.5197	1.1572	0.7214	3.3072	0.4266
2,4	LM	1.0158	0.9910	0.5221	0.9263	1.5252	0.1717
	GP	2.9163	2.3631	1.4808	0.8449	2.8834	0.5617
3,4	LM	0.7123	0.7123	0.5266	0.9124	1.5011	0.1633
	GP	2.7110	2.1946	1.7376	0.9026	2.7573	0.6300

From Table 8, among all configurations, Combination 1 achieves the best overall performance, yielding the lowest values of NLPD, RMSE, MAE, and CRPS, together with the narrowest prediction intervals (PINAW). Combination 2 shows comparable but slightly inferior performance, while Combination 3 and Combination 4 exhibit noticeably higher error metrics and wider prediction intervals.

Pairwise station combinations generally lead to higher prediction errors and greater uncertainty than single-station cases, with Combination 2–3 showing the weakest performance across most metrics.

Table 8. Performance metrics for different combinations (best in **bold**, worst in *italics*).

Combination	NLPD	PINAW	PINAD	RMSE	MAE	CRPS	CVG	Phi
1	1.895639	0.252707	0.00931	1.14983	0.87548	0.75921	0.94999	29.8
2	1.920527	0.288044	-0.02668	1.26153	1.00768	0.80326	0.94999	29.1
3	2.149783	0.451443	-0.06350	2.05730	1.58239	1.15320	0.94999	21.7
4	<i>2.454196</i>	<i>0.558282</i>	<i>0.10413</i>	<i>2.78430</i>	<i>2.33373</i>	<i>1.60778</i>	0.94999	27.35
1-2	2.319954	0.511679	-0.00712	2.44230	1.96784	1.38523	0.94999	26.05
1-3	2.398002	0.591682	-0.05786	2.65672	1.99771	1.46377	0.94999	21.4
1-4	2.442053	0.559535	0.07942	2.74557	2.24951	1.56815	0.94999	28.15
2-3	<i>2.792790</i>	<i>0.679887</i>	-0.07146	<i>3.16869</i>	<i>2.50263</i>	<i>1.83119</i>	0.94999	23.55
2-4	2.464719	0.603344	0.06321	2.87521	2.33282	1.63096	0.94999	25.95
3-4	2.429288	0.545177	0.02086	2.70387	2.19068	1.53779	0.94999	26.5

Across all configurations, the empirical coverage (CVG) remains close to the nominal level of 0.95, indicating that the probabilistic forecasts are well calibrated regardless of the station combination. The relative improvement metric Φ indicates consistent performance gains over the linear baseline, with the largest improvements observed for single-station configurations.

Figure 11 presents the split conformal prediction results for Combination 1 and Combination 1–2, while the corresponding quantitative performance metrics are reported in Table 8. For Combination 1, the predicted mean closely tracks the observed wind speed, with relatively narrow prediction intervals, consistent with the low RMSE (1.15), MAE (0.88), CRPS (0.76), and PINAW (0.25) values reported in the table. In contrast, Combination 1–2 exhibits wider prediction intervals and increased variability in the predicted mean, which aligns with the higher RMSE (2.44), MAE (1.97), CRPS (1.39), and PINAW (0.51) observed in Table 8. In both cases, the empirical coverage shown in Figure 11 is consistent with the nominal level, in agreement with the CVG value of approximately 0.95 reported in the table. The split conformal prediction method provides prediction intervals (shown as shaded regions) that effectively capture the variability in the observed wind speed data (black line). The predicted values (red dashed line) closely follow the actual measurements, and the uncertainty bands adapt over time, widening during more volatile periods and narrowing during more stable periods. This shows that split conformal prediction generates well-calibrated, dynamic-width intervals that reflect the true uncertainty in predictions across disparate time indices.

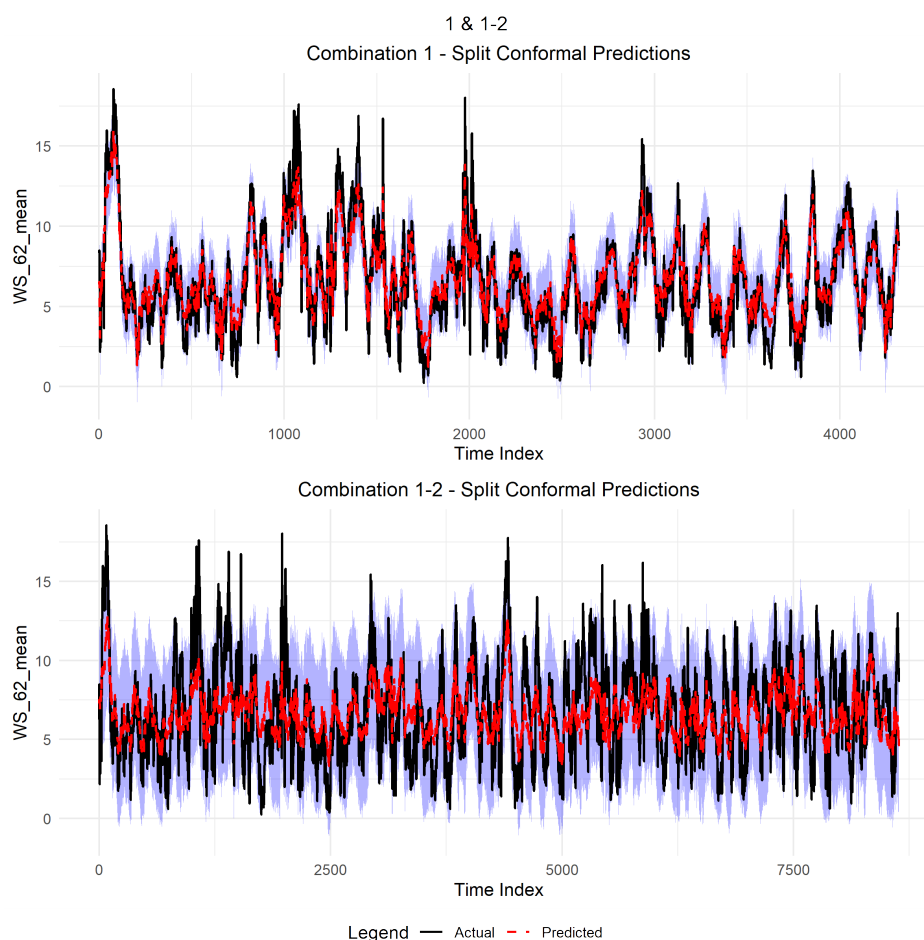


Figure 11. Prediction plots for cluster 2 high performing stations.

3.2.2. Results for the Hybrid Linear-Gaussian Process

Clusters Setup

The dataset used for the hybrid model was divided into two clusters based on elevation, longitude, and latitude to account for spatial heterogeneity in wind behaviour. Cluster 1 consists of lower to mid-elevation stations (152–848m), while Cluster 2 contains higher-elevation stations (770–1806m). This is important for accounting for heterogeneity in wind speed dynamics. Model performance was evaluated separately for each cluster using multiple validation splits. This strategy allows assessment of whether the proposed hybrid framework performs consistently across regimes with different statistical and physical characteristics.

Validation Setup

The proposed hybrid modelling framework was evaluated using a two-fold validation strategy. A leave-one-station-out cross-validation approach was used to evaluate model performance. In each iteration, data from one wind station were withheld for validation, while the remaining stations were used for model training. This procedure was repeated until each station had served as the validation site once. For cluster 1 (total sample size 25914), each validation set comprised 16.7% of the available observations (4319 samples per split) and for cluster 2 (total sample size 17276), each validation set comprised 25% of the available observations (4319 samples per split). For each split, a two-stage process was used. First, a linear regression model was estimated to account for large-scale spatiotemporal effects, followed by a GP analysis of the residuals. A variety of covariance structures were rigorously evaluated within the second stage, including several Matérn covariance functions with smoothness $u = 0.5, 1.5, \text{ and } 2.5$, an exponential covariance, and a Gaussian covariance. The best-performing covariance structure is based on validation performance.

Results for Cluster 1: Weak or Unstable Residual Structure

Predictive Performance

In respect of Cluster 1, the improved linear-Gaussian Process (GP) hybrid model performance was reported as mixed, with varying levels of improvement depending on the residual errors it could model relative to the purely linear model. From the results presented in Table 9, the three distinct validation set RMSEs were 3.127, 21.779, and 2.828, while the linear model reported values were 3.125, 21.431, and 2.844. Here, the enhancements reported ranged from -1.622 to 0.576, indicating that this approach is indeed useful for modelling the non-linear residuals. However, the overall improvement is limited by the varying spatial residuals in Cluster 1. In one of the folds, it reports an improvement of 0.576, despite the overall spatial correlation in the residuals.

Covariance Selection and Residual Structure

In all three model validation cases, the Gaussian (squared exponential) covariance was selected as optimal. This suggests that the residuals of Cluster 1 could be best modelled as having a stationary covariance structure. Additionally, the negative improvement in the GP fold indicates that the residuals' spatial structure is neither strong nor stationary, resulting in little improvement from the GP method. Covariance-based feature selection cannot always adequately compensate for residual heterogeneity.

Probabilistic Forecast Quality

The hybrid method provided well-calibrated uncertainty estimates. The coverage (CVG) was kept at 0.95 for all the folds, and the PINAW values ranged from 0.465 to 2.511. Low PINAD values across most folds indicate little bias in interval centring. This shows that even when point prediction was not significantly improved, the GP component still provided good predictions; hence, decisions could be made accordingly.

Implications for Wind Speed Forecasting

The results for Cluster 1 should be kept in mind for the practical application of the wind forecast method. In particular, the usefulness of the hybrid GP modelling method depends on the spatial coherence of the residuals; a lack of heterogeneity in the residuals could prevent efficient RMSE improvement within the clusters. The selection of covariance should be based on validation performance, but this may require non-stationary and/or multi-scale kernels for more accurate terrain-effect modelling. Probabilistic forecasts remain robust, and a small improvement in point predictions underscores the scope for risk-informed decision-making with hybrid models, especially in wind energy applications. The findings from cluster 1 highlight that modelling strategies should vary across clusters, accounting for trends and residual spatial structures, especially in complex terrain and heterogeneous winds.

Results for Cluster 2: Strong Spatio-Temporal Dependence

Predictive Performance

In both sets of validation data, the proposed hybrid linear-GP model significantly performed better than the linear-only model, validating the presence of underlying spatiotemporal correlation in the residual wind speed field (Table 9). For the first validation set, the proposed model showed improved performance, achieving an RMSE of 2.51, compared with the linear model, which achieved an RMSE of 2.99. The performance improvements of the proposed model were hence 16.3%. On the second set, the improvements were slightly lower, at 3.2%, as indicated by the reduced RMSE from 2.39 to 2.31.

Covariance Selection and Residual Structure

Different optimal covariance structures were selected based on the validation splits reported by the machine-learning algorithm, using the automatic covariance structure selection method implemented

in the code. In the first validation set, the Matérn covariance structure with $u = 0.5$ was consistently selected, while in the second set, the Gaussian covariance structure was preferred. The Matérn $u = 0.5$ covariance structure essentially corresponds to a rough, even non-differentiable, stochastic process, which is typically preferred for modelling abrupt changes in the underlying dynamics, for instance, those characterising localised variability in the wind field. The optimal structure found for the first set suggests that the dynamics of the residual wind field were characterised by particularly rough variability, potentially linked to factors such as gustiness and rapid regime changes. On the contrary, the optimal choice of the Gaussian kernel was observed in the second set, where the proportion of the original variability explained by the linear model was also particularly large.

Probabilistic Forecast Quality

Aside from point prediction accuracy, the hybrid models also demonstrated excellent probabilistic performance. Also, the Matérn $u = 0.5$ performed best in terms of the smallest negative log predictive density, well-centred prediction intervals, and almost nominal empirical coverage of about 95%. Even though the prediction intervals generated by the Gaussian kernel were narrower and the CRPS was slightly reduced, the model also showed some underdispersion, indicating overconfidence in the predictions. These findings reveal a balance between point prediction accuracy and probabilistic calibration, with rougher covariance structures yielding a more realistic quantification of uncertainty in wind speed predictions.

Implications for Wind Speed Forecasting

The results show the capability of hybrid linear-GP approaches with adaptive covariance selection in spatio-temporal prediction of wind speed while being physically consistent. The linear model demonstrates its capability to capture larger trends, while the GP model captures the localised stochastic component. The tendency toward rough covariance structures supports the non-smooth nature of the underlying wind-speed processes, justifying the use of Matérn kernels in this context. The presented model offers greater prediction accuracy, more precise quantification of uncertainty, and greater interpretability of predictions than traditional linear models.

Table 9. Model Performance Metrics by Clusters and Validations.

Metric	C1 Val1	C1 Val2	C1 Val3	C2 Val1	C2 Val2
Covariance Type	gaussian	gaussian	gaussian	Matern_0.5	gaussian
NLPD	2.949	185.534	11.375	2.402	2.474
PINAW	0.510	2.511	0.465	0.539	0.490
PINAD	0.048	-1.053	-0.017	0.020	-0.019
RMSE	3.127	21.779	2.828	2.506	2.309
MAE	2.649	21.403	2.293	2.015	1.866
CRPS	1.931	20.763	2.000	1.429	1.345
CVG	0.950	0.950	0.950	0.950	0.950
Linear RMSE	3.125	21.431	2.844	2.994	2.386
GP Improvement (%)	-0.065	-1.622	0.576	16.279	3.242

As shown in Figure 12, the figure compares the predictive performance of the Linear and Hybrid GP regression models across two data clusters and various validation sets. The left-hand side of the figure illustrates a representation of the Root Mean Square Error (RMSE) for the Linear Model and the Hybrid GP model. As is apparent, the RMSE is a vital measure that effectively quantifies error. As such, lower error values are desirable compared to higher ones. This figure effectively compares the relative performance of including the GP model in improving overall wind speed prediction across different data clusters. On the other hand, the right-hand side of Figure 12 demonstrates the percentage improvement of the Hybrid GP Model over the basic Linear Model. The inclusion of the GP Model

would be beneficial if the value is positive. In contrast, a negative value would indicate cases where the GP Model does not effectively improve prediction accuracy.

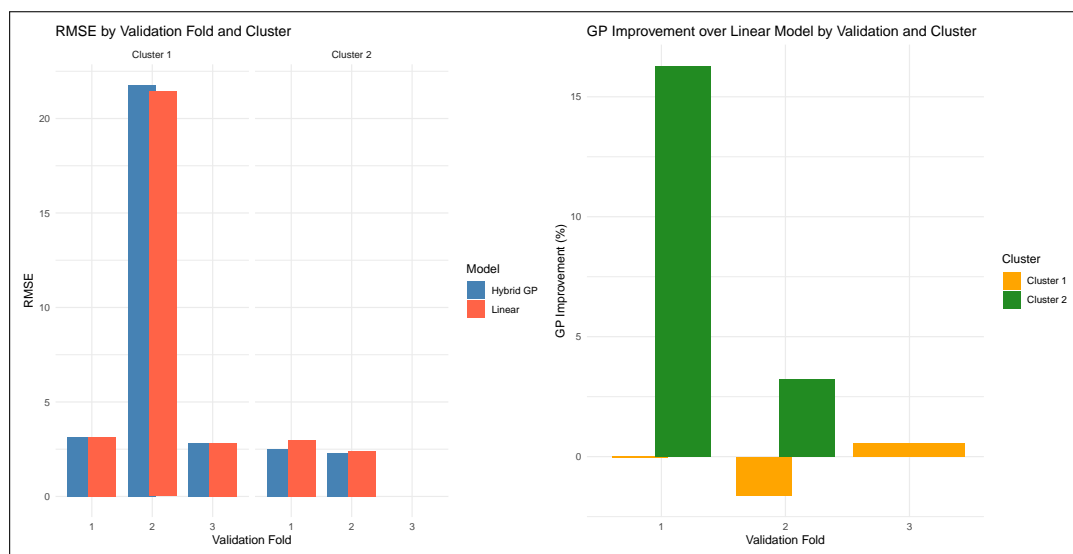


Figure 12. Comparison of predictive performance between Linear and Hybrid GP models across clusters and validation sets.

From the appendix, figure A1 presents the time series plot for the independent variable windspeed. Figures A2-A5 shows the results of the hybrid model from both clusters with the covariances selected as well as the improvement from the linear model.

4. Discussion

4.1. Discussion of Results

The proposed hybrid framework was tested using two different clusters of wind speed. Although the improvements were significant in the second cluster, the results of the first cluster show very different behaviour. In fact, the hybrid framework improved over the linear baseline in only one of the three validation splits, with an average change in RMSE of -0.4% . In one of the validation splits, the GP's performance significantly deteriorated, and the RMSE exceeded 20. The selection procedure reliably identified the Gaussian kernel as the optimal one for this cluster. This confirms our previous assumption that the variability present in this cluster was mostly smooth and weakly correlated. On the other hand, marginal improvements and even degradations in performance suggest that this wind speed variability lacks sufficient spatio-temporal structure for a GP-based model. Thus, we have successfully demonstrated the limitations of employing GP-based residual modelling techniques and shown how this particular framework can identify them through this covariance selection and validation structure.

From an operational perspective, these results reinforce the relevance of regime-aware modelling approaches. Excessive and indiscriminate use of complex stochastic models may incur unnecessary computational cost and degrade forecast quality.

4.2. Limitations and Future Work

Firstly, the stationarity of the process within clusters, as assumed by the Gaussian Process Residual model, may not account for sudden regime changes characteristic of the atmosphere. This issue may be particularly concerning for Cluster 1, as suggested by the validation results. Secondly, the Gaussian Process (GP) method may be sensitive to sparsity/outliers, as evidenced by occasional decreases in model performance as measured by validation. However, it should be noted that this issue is mitigated by adaptive covariance selection. Lastly, although the two-stage method may be more interpretable and efficient, it does not use any physics-based constraints.

Future work may involve generalising the method to non-stationary and non-separable covariance functions, using anisotropic functions that depend on wind direction, and developing approximations of the GP method for operational use.

5. Conclusions

We had two fundamentally different wind regimes obtained by clustering the wind stations. For cluster 1, the linear model is already near-optimal, Residuals are weak, noisy, or unstable and the GP often degrades performance. Smooth kernels selected, but were not very helpful. For Cluster 2 (strong spatio-temporal dependence, rough residuals), GP provides large gains, and Matérn kernels dominate. This leads to the conclusion that the effectiveness of Gaussian process residual modelling is regime-dependent and varies across wind speed clusters.

This study applied a hybrid linear–Gaussian Process model to daily wind speed data from South African stations, clustered by elevation and geographic location. Key conclusions are: Cluster 1 (low-to-mid elevation) showed mixed GP improvements due to heterogeneous residuals, but probabilistic forecasts were reliable. Cluster 2 (high elevation) experienced consistent improvements, with GP capturing coherent residual patterns and enhancing both point and probabilistic predictions. Covariance selection is critical; Gaussian and Matérn 0.5 kernels were optimal depending on residual structure. The hybrid model effectively quantifies forecast uncertainty, making it valuable for risk-aware wind energy management.

Future work should explore non-stationary or multi-scale kernels, as well as additional topographic or meteorological features, to further improve predictive performance in heterogeneous terrain.

Author Contributions: Conceptualisation, THT and CS.; methodology, THT.; software, THT.; validation, THT, CS, TBD and SN; formal analysis, THT; investigation, THT, CS, TBD and SN.; data curation, THT; writing—original draft preparation, THT; writing—review and editing, THT, CS, TBD and SN; visualisation, THT.; supervision, CS, TBD and SN; project administration, CS, TBD and SN. All authors have read and agreed to the published version of this manuscript.

Funding: This research received no external funding.

Data Availability Statement: The data in brief can be accessed at <https://github.com/csigauke/A-Hybrid-Linear--Gaussian-Process-Framework-for-Spatio-Temporal-Wind-Speed-Forecasting>.

Acknowledgments: The authors thank the anonymous reviewers for their helpful comments on this paper.

Conflicts of Interest: The authors declare no conflicts of interest

Abbreviations

The following abbreviations are used in this manuscript:

MAE	Mean Absolute Error
MASE	Mean Absolute Scaled Error
RMSE	Root Mean Squared Error
MCMC	Markov Chain Monte Carlo
CRPS	Continuous Ranked Probability Score
PICP	probability of prediction interval coverage
NLPD	Negative Log Predictive Density
PINAW	normalised average width of prediction interval
PINAD	normalised average deviation of prediction interval

Appendix A. Supplementary Plots

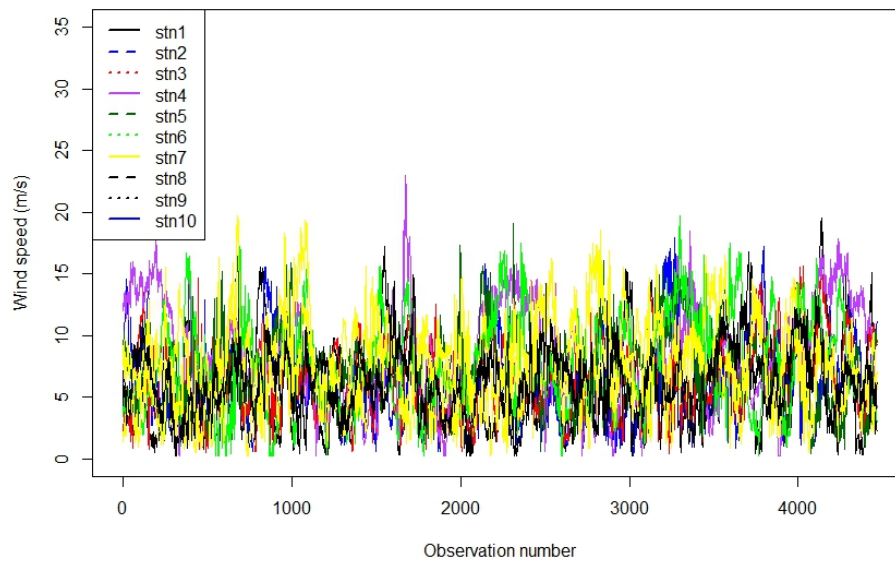


Figure A1. Time series plot for all the stations.

Hybrid model plots

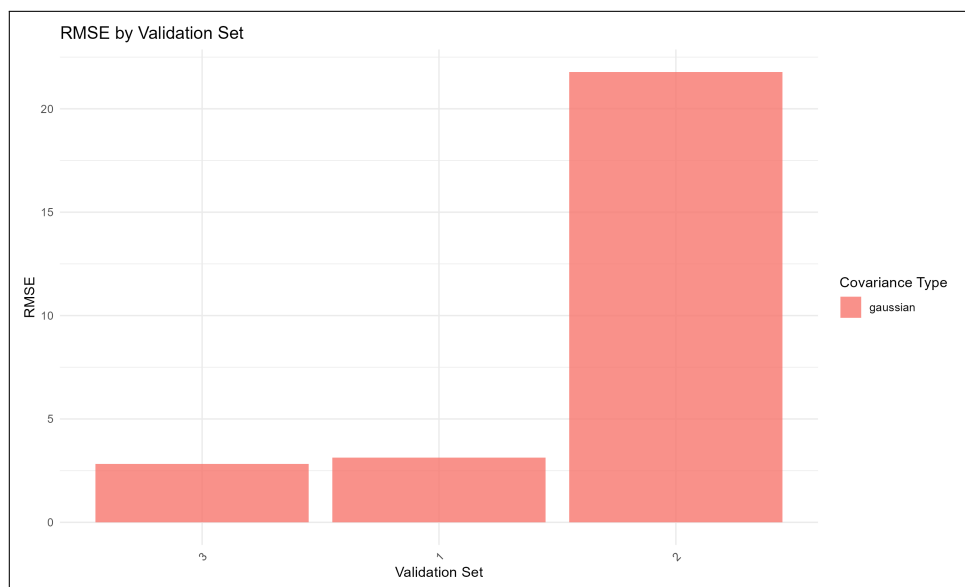


Figure A2. RMSE for cluster 1

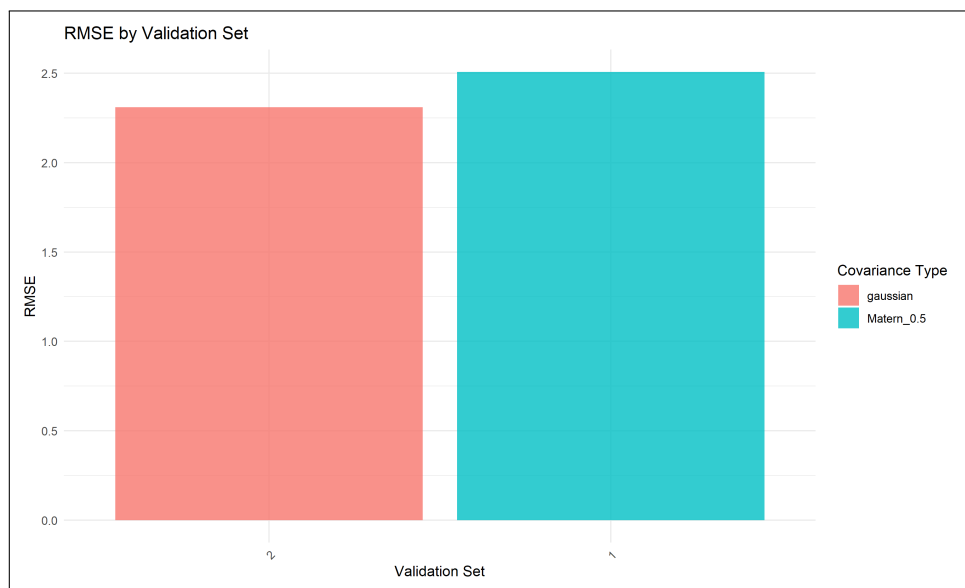


Figure A3. RMSE for cluster 2

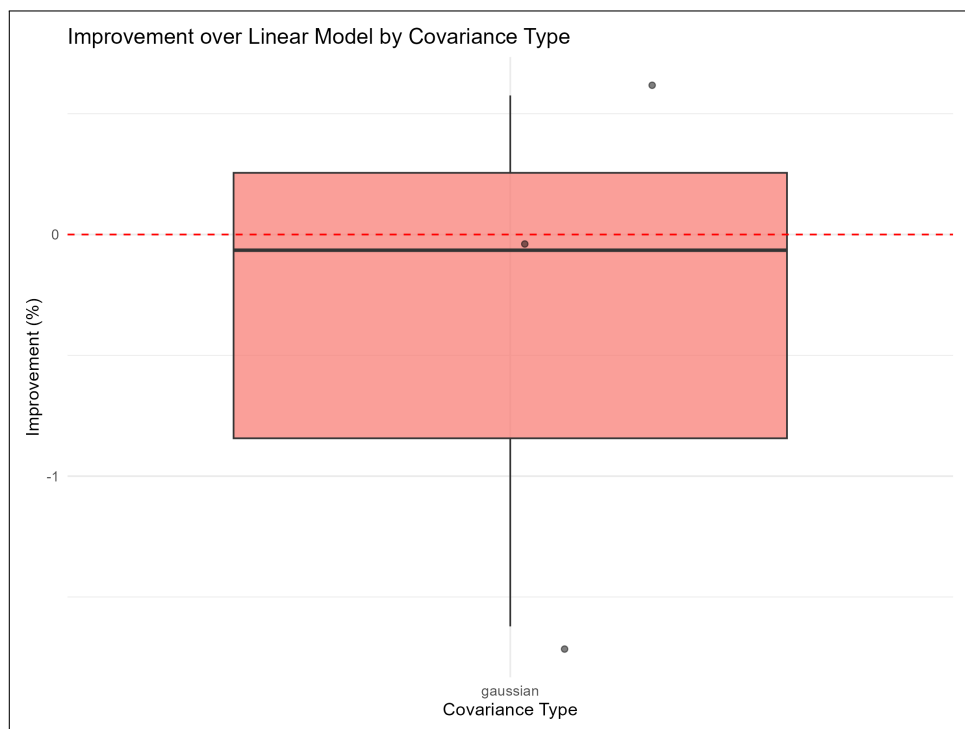


Figure A4. Improvement for cluster 1

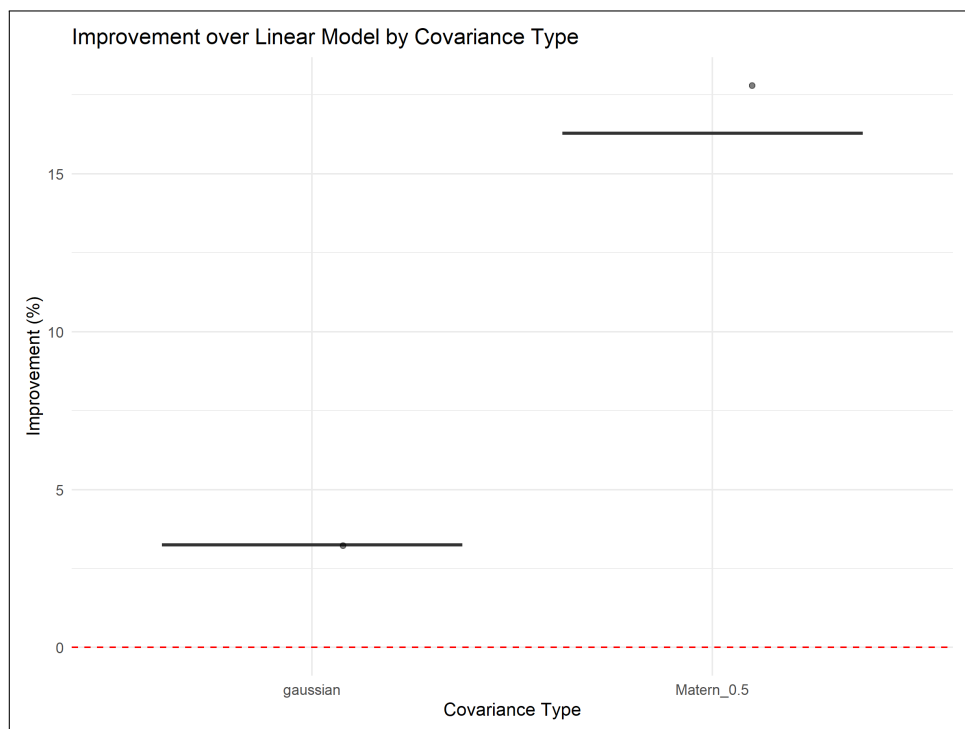


Figure A5. Improvement for cluster 2

Appendix B. Modelling specification-An Extensive Hierarchical Bayesian Formulation with MCMC Sampling.

Appendix B.1. Hybrid Model Specification

Let $y(\mathbf{s}, t)$ denote the response variable wind speed at location $\mathbf{s} \in \mathcal{S} \subset \mathbb{R}^2$ and time $t \in \mathcal{T} \subset \mathbb{R}$. The complete model is:

$$y(\mathbf{s}, t) = \underbrace{\mathbf{x}(\mathbf{s}, t)^\top \boldsymbol{\beta}}_{\text{Linear trend}} + \underbrace{f(\mathbf{s}, t)}_{\text{GP component}} + \epsilon(\mathbf{s}, t) \quad (\text{A1})$$

where:

- $\mathbf{x}(\mathbf{s}, t) \in \mathbb{R}^p$ is a vector of covariates
- $\boldsymbol{\beta} \in \mathbb{R}^p$ are regression coefficients
- $f(\mathbf{s}, t) \sim GP(0, K(\cdot, \cdot; \boldsymbol{\theta}))$ is a zero-mean Gaussian process
- $\epsilon(\mathbf{s}, t) \sim N(0, \tau^2)$ is the measurement error (nugget effect)
- $K(\cdot, \cdot; \boldsymbol{\theta})$ is the covariance function with parameters $\boldsymbol{\theta}$

Appendix B.2. Two-Stage Estimation Procedure

Appendix B.2.1. Stage 1: Linear Regression

The very first attempt at modelling a spatio-temporal response variable is to consider the linear regression models. First, estimate the linear component via ordinary least squares:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \quad (\text{A2})$$

The linear predictions are:

$$\hat{y}_L(\mathbf{s}, t) = \mathbf{x}(\mathbf{s}, t)^\top \hat{\boldsymbol{\beta}} \quad (\text{A3})$$

Residuals for GP modelling:

$$r(\mathbf{s}, t) = y(\mathbf{s}, t) - \hat{y}_L(\mathbf{s}, t) \quad (\text{A4})$$

Appendix B.2.2. Stage 2: Gaussian Process on Residuals

Model the residuals as:

$$r(\mathbf{s}, t) = f(\mathbf{s}, t) + \epsilon(\mathbf{s}, t) \quad (\text{A5})$$

where $f(\mathbf{s}, t) \sim GP(0, K(\cdot, \cdot; \boldsymbol{\theta}))$.

Appendix B.3. Covariance Structures

Appendix B.3.1. General Form of Covariance Functions

For spatio-temporal modelling, we consider separable covariance functions:

$$\text{Cov}(f(\mathbf{s}, t), f(\mathbf{s}', t')) = \sigma^2 \cdot C_s(\mathbf{s}, \mathbf{s}'; \boldsymbol{\theta}_s) \cdot C_t(t, t'; \boldsymbol{\theta}_t) \quad (\text{A6})$$

where σ^2 is the marginal variance.

Appendix B.3.2. Matérn Family of Covariance Functions

The **Matern covariance function** is defined as:

$$K_{\text{Matern}}(h; \nu, \phi) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}h}{\phi} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}h}{\phi} \right) \quad (\text{A7})$$

where:

- $h = \|\mathbf{s} - \mathbf{s}'\|$ is the Euclidean distance
- $\nu > 0$ is the smoothness parameter
- $\phi > 0$ is the range parameter
- $\sigma^2 > 0$ is the variance parameter
- $K_\nu(\cdot)$ is the modified Bessel function of the second kind

Appendix B.3.3. Special Cases

1. Exponential covariance ($\nu = 0.5$):

$$K_{\text{Exp}}(h; \phi) = \sigma^2 \exp\left(-\frac{h}{\phi}\right) \quad (\text{A8})$$

2. Matern with $\nu = 1.5$:

$$K_{\nu=1.5}(h; \phi) = \sigma^2 \left(1 + \frac{\sqrt{3}h}{\phi}\right) \exp\left(-\frac{\sqrt{3}h}{\phi}\right) \quad (\text{A9})$$

3. Matern with $\nu = 2.5$:

$$K_{\nu=2.5}(h; \phi) = \sigma^2 \left(1 + \frac{\sqrt{5}h}{\phi} + \frac{5h^2}{3\phi^2}\right) \exp\left(-\frac{\sqrt{5}h}{\phi}\right) \quad (\text{A10})$$

4. Gaussian/Squared Exponential ($\nu \rightarrow \infty$):

$$K_{\text{Gaussian}}(h; \phi) = \sigma^2 \exp\left(-\frac{h^2}{2\phi^2}\right) \quad (\text{A11})$$

- **Exponential covariance** ($\nu = 0.5$) models rough processes with less smoothness.
- **Gaussian covariance** ($\nu \rightarrow \infty$) assumes very smooth processes.
- **Matern covariance** with intermediate ν balances smoothness and flexibility.

Appendix B.3.4. Rational Quadratic Covariance

$$K_{\text{RQ}}(h; \alpha, \phi) = \sigma^2 \left(1 + \frac{h^2}{2\alpha\phi^2} \right)^{-\alpha} \quad (\text{A12})$$

where $\alpha > 0$ controls the smoothness.

Appendix B.4. Likelihood and Estimation

Appendix B.4.1. Gaussian Process Likelihood

For n observations $\mathbf{r} = (r_1, \dots, r_n)^\top$, the log-likelihood is:

$$\log p(\mathbf{r} | \boldsymbol{\theta}) = -\frac{1}{2} \left[n \log(2\pi) + \log |\mathbf{K}_\theta + \tau^2 \mathbf{I}| + \mathbf{r}^\top (\mathbf{K}_\theta + \tau^2 \mathbf{I})^{-1} \mathbf{r} \right] \quad (\text{A13})$$

where \mathbf{K}_θ is the $n \times n$ covariance matrix with entries $K_\theta(\mathbf{s}_i, \mathbf{s}_j)$.

Appendix B.4.2. Maximum Likelihood Estimation

The MLE estimates are obtained by:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{r} | \boldsymbol{\theta}) \quad (\text{A14})$$

Appendix B.4.3. Predictive Distribution

For a new location \mathbf{s}_0 , the predictive distribution is Gaussian:

$$p(r(\mathbf{s}_0) | \mathbf{r}) = N(\mu_0, \sigma_0^2) \quad (\text{A15})$$

with:

$$\mu_0 = \mathbf{k}_0^\top (\mathbf{K} + \tau^2 \mathbf{I})^{-1} \mathbf{r} \quad (\text{A16})$$

$$\sigma_0^2 = K(\mathbf{s}_0, \mathbf{s}_0) - \mathbf{k}_0^\top (\mathbf{K} + \tau^2 \mathbf{I})^{-1} \mathbf{k}_0 + \tau^2 \quad (\text{A17})$$

where $\mathbf{k}_0 = [K(\mathbf{s}_0, \mathbf{s}_1), \dots, K(\mathbf{s}_0, \mathbf{s}_n)]^\top$.

Appendix B.5. Hierarchical Bayesian Formulation

Appendix B.5.1. Complete Bayesian Model

$$\text{Likelihood: } \mathbf{y} | \boldsymbol{\beta}, f, \tau^2 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{f}, \tau^2 \mathbf{I}) \quad (\text{A18})$$

$$\text{GP Prior: } \mathbf{f} | \sigma^2, \phi, \nu \sim N(\mathbf{0}, \mathbf{K}(\sigma^2, \phi, \nu)) \quad (\text{A19})$$

$$\text{Priors: } \boldsymbol{\beta} \sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}) \quad (\text{A20})$$

$$\tau^2 \sim \text{IG}(a_\tau, b_\tau) \quad (\text{A21})$$

$$\sigma^2 \sim \text{IG}(a_\sigma, b_\sigma) \quad (\text{A22})$$

$$\phi \sim \text{LN}(\mu_\phi, \sigma_\phi^2) \quad (\text{A23})$$

$$\nu \sim \text{Gamma}(a_\nu, b_\nu) \quad (\text{A24})$$

Appendix B.5.2. Posterior Distribution

$$p(\boldsymbol{\beta}, \mathbf{f}, \boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{f}, \tau^2) \cdot p(\mathbf{f} | \boldsymbol{\theta}) \cdot p(\boldsymbol{\beta}) \cdot p(\boldsymbol{\theta}) \quad (\text{A25})$$

Appendix B.5.3. Gibbs Sampling

1. **Sample β** (conjugate):

$$\beta | \cdot \sim N(\mu_\beta, \Sigma_\beta) \quad (\text{A26})$$

$$\Sigma_\beta = \left(\frac{1}{\sigma_\beta^2} \mathbf{I} + \frac{1}{\tau^2} \mathbf{X}^\top \mathbf{X} \right)^{-1} \quad (\text{A27})$$

$$\mu_\beta = \frac{1}{\tau^2} \Sigma_\beta \mathbf{X}^\top (\mathbf{y} - \mathbf{f}) \quad (\text{A28})$$

2. **Sample \mathbf{f}** (conjugate):

$$\mathbf{f} | \cdot \sim N(\mu_f, \Sigma_f) \quad (\text{A29})$$

$$\Sigma_f = \left(\mathbf{K}^{-1} + \frac{1}{\tau^2} \mathbf{I} \right)^{-1} \quad (\text{A30})$$

$$\mu_f = \frac{1}{\tau^2} \Sigma_f (\mathbf{y} - \mathbf{X}\beta) \quad (\text{A31})$$

3. **Sample covariance parameters** (Metropolis-Hastings):

$$\theta^{(t+1)} = \theta^{(t)} + \epsilon, \quad \epsilon \sim N(0, \sigma_{\text{mh}}^2) \quad (\text{A32})$$

Accept with probability:

$$\alpha = \min \left(1, \frac{p(\mathbf{y} | \theta^*) p(\theta^*)}{p(\mathbf{y} | \theta^{(t)}) p(\theta^{(t)})} \right) \quad (\text{A33})$$

Appendix B.6. Model Comparison and Selection

Appendix B.6.1. Information Criteria

$$\text{AIC} = 2k - 2 \log \mathcal{L} \quad (\text{A34})$$

$$\text{BIC} = k \log n - 2 \log \mathcal{L} \quad (\text{A35})$$

$$\text{WAIC} = -2 \left(\sum_{i=1}^n \log E[p(y_i | \boldsymbol{\theta})] - \sum_{i=1}^n \text{Var}[\log p(y_i | \boldsymbol{\theta})] \right) \quad (\text{A36})$$

Appendix B.6.2. Cross-Validation

For K-fold CV:

$$\text{CV} = \frac{1}{n} \sum_{k=1}^K \sum_{i \in \mathcal{D}_k} (y_i - \hat{y}_{-k,i})^2 \quad (\text{A37})$$

Appendix B.7. Uncertainty Quantification

Appendix B.7.1. Prediction Intervals

For a new location \mathbf{s}_0 , the predictive distribution of the response is Gaussian. The $(1 - \alpha)$ prediction interval:

$$\hat{y}(\mathbf{s}_0) \pm z_{1-\alpha/2} \sqrt{\text{Var}[\hat{y}(\mathbf{s}_0)]} \quad (\text{A38})$$

where $\text{Var}[\hat{y}(\mathbf{s}_0)] = \sigma_0^2 + \tau^2$. This provides both a **point prediction** (μ_0) and an **uncertainty measure** (σ_0^2) for the new location.

Appendix B.8. Conformal Prediction

1. Compute residuals: $e_i = |y_i - \hat{y}_i|$, $i \in \mathcal{D}_{\text{cal}}$

2. Find α -quantile: $q_{1-\alpha} = \text{Quantile}(1 - \alpha; \{e_i\})$
3. Prediction interval: $\text{PI}(\mathbf{s}_0) = [\hat{y}(\mathbf{s}_0) - q_{1-\alpha}, \hat{y}(\mathbf{s}_0) + q_{1-\alpha}]$

Appendix B.9. Computational Aspects

Appendix B.9.1. Scalability Challenges

Exact GP: $\mathcal{O}(n^3)$. Solutions:

- **Low-rank approximations:** $\mathbf{K} \approx \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{nm}^\top$, $\mathcal{O}(nm^2)$
- **Sparse covariance methods**
- **Stochastic Variational Inference**

Appendix B.9.2. Numerical Stability

- Add jitter: $\mathbf{K} + \delta \mathbf{I}$
- Use Cholesky decomposition with pivoting
- Compute log-determinant via Cholesky factors

Appendix B.10. Theoretical Properties

Appendix B.10.1. Consistency

Under regularity conditions, as $n \rightarrow \infty$:

$$\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0 \quad (\text{A39})$$

$$\hat{y}(\mathbf{s}) \xrightarrow{p} E[y(\mathbf{s}) \mid \mathcal{D}] \quad (\text{A40})$$

$$\text{PI}_{1-\alpha}(\mathbf{s}) \xrightarrow{p} \text{True PI}_{1-\alpha}(\mathbf{s}) \quad (\text{A41})$$

Appendix B.10.2. Asymptotic Normality

For fixed-domain asymptotics:

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N(\mathbf{0}, \mathcal{I}(\boldsymbol{\theta}_0)^{-1}) \quad (\text{A42})$$

where $\mathcal{I}(\boldsymbol{\theta}_0)$ is the Fisher information matrix.

References

1. M. Yang, Y. Jiang, J. Che, Z. Han, and Q. Lv, *Short-Term Forecasting of Wind Power Based on Error Traceability and Numerical Weather Prediction Wind Speed Correction*, *Electronics*, vol. 13, no. 8, 1559, 2024, <https://doi.org/10.3390/electronics13081559>
2. K. Szostek, D. Mazur, G. Dralus, and J. Kuszniier, *Analysis of the Effectiveness of ARIMA, SARIMA, and SVR Models in Time Series Forecasting: A Case Study of Wind Farm Energy Production*, *Energies*, vol. 17, no. 19, 4803, 2024, <https://doi.org/10.3390/en17194803>
3. J. Xu, Z. Xiao, Z. Lin, et al., *System bias correction of short-term hub-height wind forecasts using the Kalman filter*, *Protection and Control of Modern Power Systems*, vol. 6, Article 37, 2021, <https://doi.org/10.1186/s41601-021-00214-x>
4. Liu, Z.; Guo, H.; Zhang, Y.; Zuo, Z. *A Comprehensive Review of Wind Power Prediction Based on Machine Learning: Models, Applications, and Challenges*, *Energies*, vol. 18, no. 2, 350, 2025, <https://doi.org/10.3390/en18020350>
5. Che, G.; Zhou, D.; Wang, R.; Zhou, L.; Zhang, H.; Yu, S. *Wind Energy Assessment in Forested Regions Based on the Combination of WRF and LSTM-Attention Models*, *Sustainability*, vol. 16, no. 2, 898, 2024, <https://doi.org/10.3390/su16020898>
6. R. J. Bessa, V. Miranda, and J. Gama, *The Role of Predictive Distributions in Modern Wind Power Forecasting Frameworks*, *International Journal of Forecasting*, vol. 38, no. 3, pp. 601–615, 2022. <https://doi.org/10.1016/j.ijforecast.2021.11.007>

7. N. Elshaboury and H. Elmousalami, *Wind speed and power forecasting using Bayesian optimised machine learning models in Gabal Al-Zayt, Egypt*, Scientific Reports, vol. 15, Article 28500, 2025, <https://doi.org/10.1038/s41598-025-13140-x>
8. Yang, Q.; Huang, G.; Li, T.; Xu, Y.; Pan, J. A Novel Short-Term Wind Speed Prediction Method Based on Hybrid Statistical-Artificial Intelligence Model with Empirical Wavelet Transform and Hyperparameter Optimization. *Journal of Wind Engineering and Industrial Aerodynamics* 2023, 239, 105499. <https://doi.org/10.1016/j.jweia.2023.105499>
9. H. Cai, X. Jia, J. Feng, W. Li, Y.-M. Hsu, and J. Lee, *Gaussian Process Regression for Numerical Wind Speed Prediction Enhancement*, Renewable Energy, vol. 146, pp. 2112–2123, 2020. <https://doi.org/10.1016/j.renene.2019.08.018>
10. M. J. Heaton, A. E. Gelfand, D. J. Vecchia, and J. R. Guinness, *A Case Study Competition Among Methods for Analysing Large Spatial Data*, Journal of Agricultural, Biological and Environmental Statistics, vol. 24, no. 3, pp. 398–425, 2019. <https://doi.org/10.1007/s13253-018-00348-w>
11. B. Wang, L. Yan, Q. Rong, J. Chen, P. Shen, and X. Duan, *Dynamic Gaussian Process Regression for Spatio-Temporal Data Based on Local Clustering*, Chinese Journal of Aeronautics, vol. 37, no. 12, pp. 245–257, 2024. <https://doi.org/10.1016/j.cja.2024.06.026>
12. C. Sigauke, E. Chandiwana, and A. Bere, *Spatio-Temporal Forecasting of Global Horizontal Irradiance Using Bayesian Inference*, Applied Sciences, vol. 13, no. 1, 201, 2023. <https://doi.org/10.3390/app13010201>
13. Chen, Y.; et al. *Hybrid Forecasting Method for Wind Power Integrating Spatial Correlation and Corrected Numerical Weather Prediction*, Applied Energy, vol. 293, 116951, 2021. <https://doi.org/10.1016/j.apenergy.2021.116951>
14. Smith, J.; et al. *Probabilistic Wind Power Forecasting for Newly-Built Wind Farms Based on Multi-Task Gaussian Process Method*, Renewable Energy, vol. 217, 119054, 2023. <https://doi.org/10.1016/j.renene.2023.119054>
15. S. Dimitrov, Q. Li, and M. Sanchez, *Toward Integrated Spatio-Temporal Wind Forecasting: Addressing Turbine Wake Effects*, Wind Energy Science, vol. 8, no. 1, pp. 101–117, 2023. <https://doi.org/10.5194/wes-8-101-2023>
16. Y. Wang, Q. Hu, and Z. Meng, *Deep Learning-Based Wind Speed Forecasting with Temporal Feature Extraction*, Renewable Energy, vol. 189, pp. 731–744, 2022. <https://doi.org/10.1016/j.renene.2022.03.032>
17. CSIR, *WASA Data Portal*, <https://wasadata.csir.co.za/wasa1/WASADData>, accessed 2025.
18. B. Hopkins and J. G. Skellam, *A new method for determining the type of distribution of plant individuals*, Annals of Botany, vol. 18, no. 2, pp. 213–227, 1954. <https://doi.org/10.1093/oxfordjournals.aob.a083391> Available at: <https://academic.oup.com/aob/article-abstract/18/2/213/277917>
19. N. Meinshausen, *Relaxed Lasso*, Computational Statistics & Data Analysis, vol. 52, no. 1, pp. 374–393, 2007. <https://doi.org/10.1016/j.csda.2006.12.009>
20. C. K. I. Williams and C. E. Rasmussen, *Gaussian Processes for Machine Learning*, MIT Press, Cambridge, MA, 2006. <http://www.gaussianprocess.org/gpml/>
21. T. Gneiting and A. E. Raftery, *Strictly Proper Scoring Rules, Prediction, and Estimation*, Journal of the American Statistical Association, vol. 102, no. 477, pp. 359–378, 2007. <https://doi.org/10.1198/016214506000001437>

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.