

# Validation of a methodology for post-construction Energy Yield Assessment of an operational wind farm

M. Costa<sup>1</sup>, T. Rocha<sup>2</sup>, J. Mendonça<sup>3</sup>, R. Pilão<sup>4</sup> and P. Pinto<sup>1</sup>

<sup>1</sup> MEGAJOULE S.A.

Rua do Divino Salvador de Moreira, 255; 4470-105 Maia (Portugal)  
mariana.costa@megajoule.pt,  
paulo.pinto@megajoule.pt

<sup>2</sup>INESC TEC, Centre for Power and Energy Systems

Rua Dr. Roberto Frias, 4200-465 Porto (Portugal)  
tomas.rocha@inesctec.pt

<sup>3</sup> SIIS, Instituto Politécnico do Porto,

R. Dr. Roberto Frias, 4200-465 Porto (Portugal)  
jpm@isep.ipp.pt

<sup>4</sup> CIETI, ISEP Instituto Politécnico do Porto

Rua Dr. António Bernardino de Almeida, 341; 4200-072 Porto (Portugal)  
rmp@isep.ipp.pt

## Abstract.

The uncertainty associated with the prospective Energy Yield Assessment (EYA) of a wind farm may be reduced by re-estimating the energy yield after it enters normal operation. This study aims to validate a simple methodology for conducting post-construction EYA of an operational wind farm. The proposed methodology derives a linear relationship between a historical source of wind speed data and the observed wind farm production on a monthly basis. In a first stage, the impact of different data sources on the accuracy of the Long-Term energy yield estimate was assessed. Results suggest that the determination coefficient  $R^2$  is a reliable indicator for selecting the most adequate source of historical wind speed data to be used in the Long-Term energy yield estimate. In a second stage, the model was validated from a statistical point of view by testing the premises of the linear regression model, namely the significance of the linear correlation (ANOVA test), and normally-distributed (Shapiro-Wilk test), non-self-correlated (Durbin-Watson), homoscedastic (Breusch-Pagan test) residuals. Results show these premises are verified for most test cases, indicating that the model is statistically robust that the model is statistically robust for most test cases.

## Key words.

Wind energy, post-construction Energy Yield Assessment, linear regression.

## 1. Introduction

Implementing a wind farm involves several challenges, including assessing its economic viability and securing financing. However, a prospective Energy Yield Assessment (EYA) has a significant underlying uncertainty which can be attributed to uncertainty in the

local measurements and intrinsic limitations of the wind models. After construction and once the wind farm enters normal operation, reassessing the annual energy yield using operational data (post-construction EYA) allows for a significant reduction in uncertainty relative to the pre-construction study. This is key if the wind farm owner wishes to refinance or sell the project.

There are several methods for the post-construction EYA of an operational wind farm, including the historical power curve method and the index method. Both methods apply the Measure-Correlate-Predict (MCP) to a historical Long-Term data source [1]. The MCP statistical method may reduce uncertainties in a wind project by correlating observed energy yield data from the operational wind farm with a Long-Term wind database. The wind database used in the MCP method can be of various types, the most common being reanalysis. The most commonly used reanalysis series are MERRA-2 and ERA5 and the main distinction between which lies in the spatial resolution, in terms of degrees, and the height at which the wind characteristics are measured [2], [3]. This approach establishes a correlation between these two data sources for a concurrent period, which is then applied to a longer dataset of historical data to estimate the long-term energy yield, as shown in Figure 1.

This study aims to validate a simple methodology for the post-construction EYA of operational wind farms, using statistical methods that analyze the main assumptions considered.

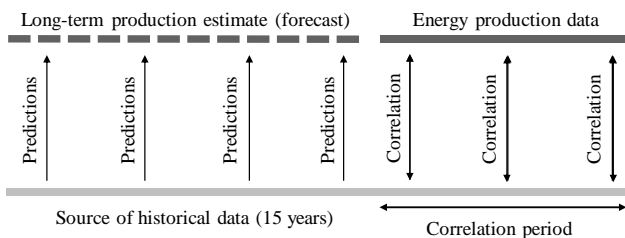


Fig. 1. Schematic demonstration of the correlation between wind farm production data and data from the historical data source, in the MCP method.

## 2. Methodology

The study database consisted of 26 wind farms with varying operational periods (4 to 20 years), distributed across different countries: Portugal, Spain, Romania, and Canada. The selection of wind farms in different parts of the world aimed to capture different wind regimes and the effect they may have on the quality of the reanalysis series itself and the estimation of long-term energy production. The general information on each wind farm used in this study is described in Table I.

Table I. Generic information for each wind farm

Country	Designation	Installed capacity [MW]	Years of operation
Portugal	P1	26	6
	P2	15	12
	P3	40	6
	P4	102	5
	P5	12	20
	P6	21.6	20
	P7	18	20
	P8	0.6	15
	P9	18.4	15
	P10	20	15
	P11	144	15
	P12	13,5	15
Spain	E1	64,7	4
	E2	100,7	4
	E3	54,4	4
	E4	128	4
	E5	45	4
	E6	20	4
	E7	49,5	4
	E8	49,5	4
	E9	30,6	4
	E10	25,2	4
	E11	30	4
Romania	R1	33,6	4
	R2	8,4	4
Canada	C1	100	9

The input data used is the actual energy production of the selected wind farms, on a monthly basis, and the availability of the wind turbines that make up the farm, for the operating period. The wind farm's actual production data is corrected for 100 % farm availability, i.e. the ratio between the monthly production value and the farm's availability in that month is calculated. This results in a series of monthly energy production data equivalent to 100 % wind farm availability.

For the central location of each wind project, 8 reanalysis series (4 of type ERA5 and 4 of type MERRA-2) were extracted for grid points closest to that central point, as shown in Figure 2. For each reanalysis data set, the observed monthly production data of each wind farm was linearly correlated with Long-Term wind data series for the simultaneous period, resulting in 8 correlations. To improve these correlations, points that are too far from the trend line or have a monthly availability value of less than 85 % are filtered out. The established correlation was applied to the selected reanalysis data series, generating a reconstructed monthly production series for the wind farm's 15 years of operation, since 15 years was considered to be representative of the Long Term. Subsequently, the annual production series was obtained by summing the monthly productions for each year, and the Long-Term energy production estimate was calculated by averaging the annual productions over the 15-year period. The method is schematized in Figure 3.

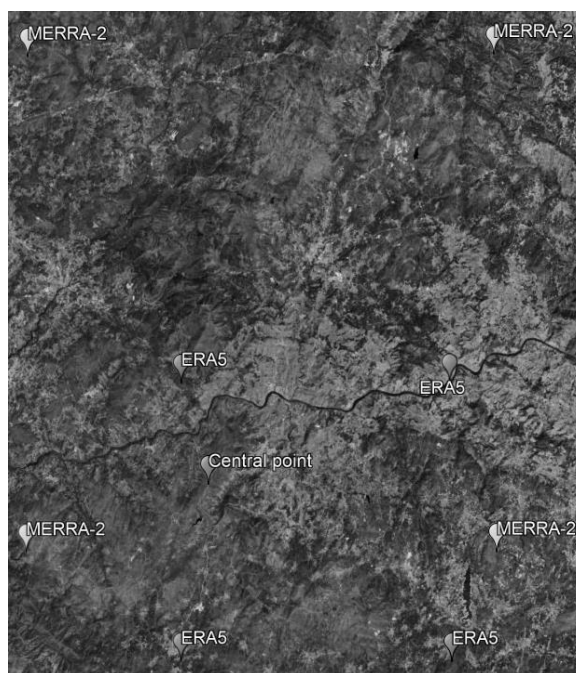


Fig. 2. Representation of the central point of a certain wind farm and the respective extracted reanalysis points.

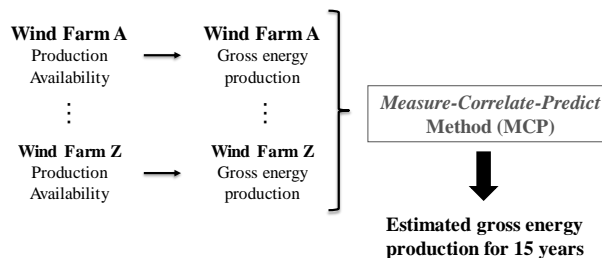


Fig. 3. MCP method schematized.

The validation of the adopted methodology was conducted by verifying the compliance with the assumptions of linear regression for the obtained correlation, through statistical analyses assessing each of the most relevant parameters. These analyses were applied to all wind farms with 4 to 15 years of operational data. The statistical tests used included:

A. *The analysis of variance (ANOVA)* is used to assess the statistical significance of the regression model, helping to determine whether the linear regression fitted to the data is statistically significant and whether the independent variables are making a significant contribution to explaining the behavior of the dependent variable (alternative hypothesis). The p-value is the statistical measure that helps determine whether there is evidence to reject the null hypothesis, so if the p-value is lower than the significance level (namely 0.05) it is possible to state that the model is statistically significant [4].

B. *The Shapiro-Wilk test* is the most powerful test for assessing the normality of a sample's distribution. If the p-value is less than the significance level (usually 0.05), the null hypothesis can be rejected, and it can be concluded that the data does not follow the theoretical distribution. In the context of this work, this test was applied to assess the normality of the residuals [5].

C. *The Durbin-Watson test* is a statistical procedure used to check for the presence of autocorrelation in the residuals of a regression model, which occurs when the residuals of a model are not independent of each other, i.e. there is a relationship between the error at one point and the error at previous points. This test calculates a statistic that varies between 0 and 4, where a value close to 2 indicates that there is no significant autocorrelation in the residuals and a value less than or greater than 2 suggests that adjacent residuals are correlated [6].

D. *The Breusch-Pagan test* is used in regression analysis to assess the variability of the model's residuals. The test assumes that the variability of the residuals in regression models is constant, calculating a statistic under the null

hypothesis of constant variability and if the p-value associated with the D-statistic is less than a certain significance level (usually 0.05), the null hypothesis can be rejected, and it can be concluded that there is evidence of non-constant variability [7].

### 3. Discussion of results

The correlations were analyzed to choose the most suitable reanalysis point for the study, with the selection criterion being the series associated with the combination of random residuals and the highest determination coefficient ( $R^2$ ). Using the wind farms with 15 years of operational data (wind farms P5 to P12) as a database, an analysis was carried out of all the reanalysis data sources extracted for these eight wind farms.

The main objective of this analysis was to investigate the relationship between the value of the  $R^2$  and the value of the percentage deviation between the estimate of Long-Term energy production obtained for each of the reanalysis series and the actual average for the wind farm (Long Term Energy Yield Error). For each of the reanalysis series extracted, the reconstructed production series was estimated by calculating the value of the Long-Term energy production estimate for this reanalysis series, and then calculating the percentage deviation between this value and the actual average value for the 15 years of operation. In this way, it is possible to relate the value of the determination coefficient of the linear relationship obtained for each reanalysis series to the deviation between the estimate of long-term energy production for that reanalysis series and the actual average for the wind farm, as shown in Figure 4.

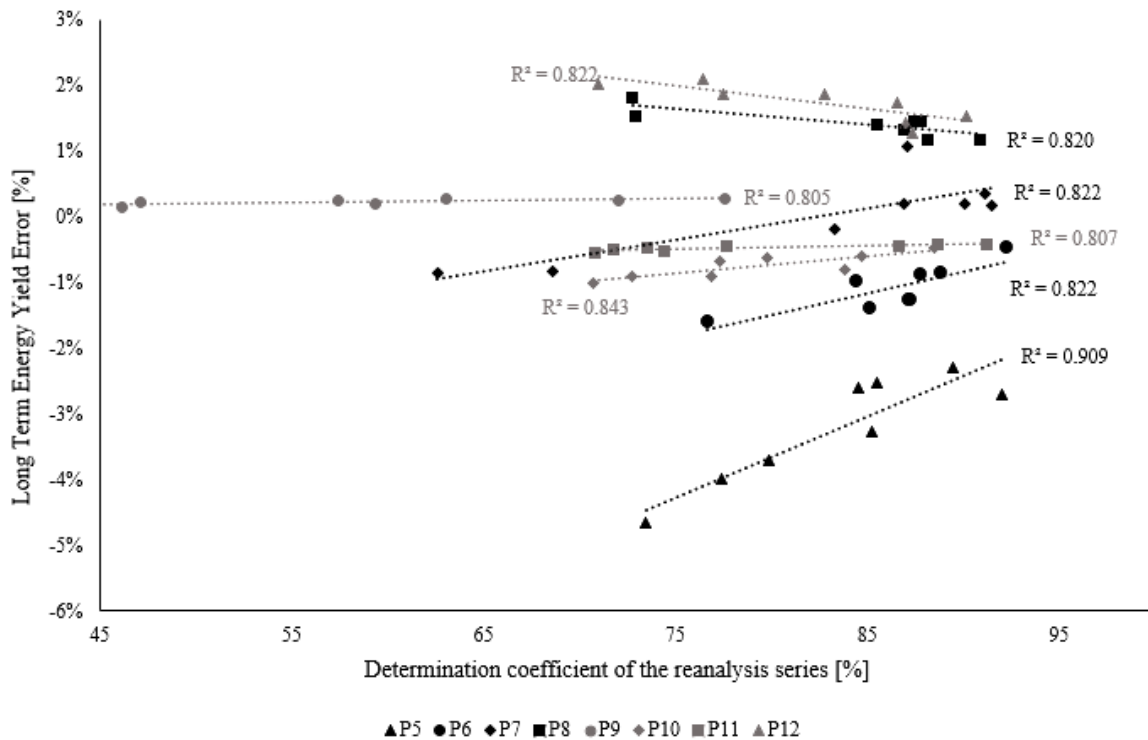


Fig. 4. Long Term Energy Yield Error, for wind farms P5 to P12.

The analysis revealed that there is a high correlation, above 80 % for the wind farms under study, between the  $R^2$  associated with each reanalysis series and Long-Term Energy Yield Error. This correlation indicates a strong relationship between these two variables. In the case of wind farms P9 and P11, there is no direct relationship as mentioned above - the Long-Term Energy Yield Error is approximately constant, regardless of the  $R^2$  value. Although the determination coefficient measures the linear relationship between two variables, real phenomena are not always strictly linear, so there may be non-linear factors that affect the estimated energy production of wind farms P9 and P11. Despite this exception, the information contained in the graphical representation in Figure 4 suggests that the  $R^2$  determination coefficient may be a good indicator for choosing the most suitable reanalysis series for studying the wind farm in question.

To verify compliance with the assumptions of the linear regression model, all the wind farms were subjected to the set of statistical tests mentioned above, which allowed us to draw relevant conclusions for the validation of the methodology used. To do this, the reanalysis series with the highest  $R^2$  for each of the wind farms under study was selected. The results obtained are shown in Table II.

Table II - Results relating to the statistical tests carried out. ANOVA: SS = statistically significant; Shapiro Wilk: ND = normal distribution and NND = non-normal distribution; Durbin Watson: AR = autocorrelation in the residuals and NAR = no autocorrelation in the residuals; Breusch Pagan: CV = constant variability in the residuals and NCV = no constant variability in the residuals

Wind Farm	ANOVA	Shapiro-Wilk	Durbin Watson	Breusch-Pagan
E1	SS	ND	NAR	CV
E2	SS	ND	NAR	CV
E3	SS	ND	NAR	CV
E4	SS	ND	NAR	NCV
E5	SS	ND	NAR	CV
E6	SS	ND	NAR	CV
E7	SS	ND	AR	CV
E8	SS	ND	NAR	CV
E9	SS	ND	NAR	CV
E10	SS	ND	NAR	CV
E11	SS	ND	NAR	CV
R1	SS	ND	NAR	NCV
R2	SS	ND	NAR	CV
C1	SS	ND	NAR	CV
P1	SS	ND	NAR	CV
P2	SS	ND	AR	CV
P3	SS	ND	NAR	CV
P4	SS	ND	NAR	NCV
P5	SS	ND	NAR	NCV
P6	SS	NND	NAR	NCV
P7	SS	ND	NAR	CV
P8	SS	ND	NAR	NCV
P9	SS	ND	NAR	NCV
P10	SS	ND	NAR	CV
P11	SS	NND	NAR	CV
P12	SS	ND	NAR	NCV

The analysis of variance (ANOVA) carried out showed that, for all parks and complexes, the linear regression model used is statistically significant. The Shapiro-Wilk test revealed that the residuals for all the parks follow a normal distribution, except for wind farms P6 and P11. According to the Durbin-Watson test, there is no autocorrelation in the residuals for most of the wind farms, except wind farms E7 and P2. Finally, the Breusch-Pagan test showed that there are eight parks (E4, R1, P4, P5, P6, P8, P9 and P16) that show evidence of non-constant variances. Although some wind farms violate the assumptions of autocorrelation in the residuals and constant variances, it can be said that, in general, the methodology can be validated.

To solve the problem of autocorrelation of the residuals, it would be necessary to assess whether there are other parameters influencing the dependent variable, in addition to the independent variable itself, and if there are, to include them in the linear regression model and perform the Durbin-Watson test again to analyze whether the problem has been solved. In the case of non-constant variances, it would be necessary to transform the equation representing the linear regression model to the logarithmic type and repeat the Breusch-Pagan test to see if the problem of non-constant variances has been eliminated. These solutions were not tested in this work because the assumptions were verified in most case studies.

#### 4. Conclusion

The first analysis suggests that  $R^2$  can serve as a reliable indicator for selecting the appropriate reanalysis series in the energy production estimation process, since the analysis carried out revealed that there is a high value of the determination coefficient, above 80 % for the wind farms under study, between the determination coefficient of the linear regression obtained for reanalysis series and the Long-Term Energy Yield Error. Overall, this indicates a strong relationship between these two variables, as the  $R^2$  value increases for a reanalysis series, the percentage deviation from the Long-Term average tends to decrease, and vice versa. There are still some cases in which there is no significant improvement in the Long-Term Energy Yield Error, so it is possible to conclude that choosing the reanalysis series with the highest  $R^2$  is a good selection principle, since it improves or maintains the accuracy of the estimate.

To verify compliance with the assumptions of the linear regression model, all wind farms underwent a set of statistical tests. The analysis of variance (ANOVA) conducted demonstrated that, for all wind farms, the linear regression model used is statistically significant. The Shapiro-Wilk test revealed that the residuals of all the wind farms follow a normal distribution, except two wind farms. According to the Durbin-Watson test, there is autocorrelation in the residuals of only two of the wind farms studied. The Breusch-Pagan test revealed evidence of non-constant variances in eight wind farms. Overall, the results obtained confirm the validity of the methodology.

## Acknowledgement

This work was supported by Multi-year financing of FCT-*Fundação para a Ciência e Tecnologia* (grant UIDB/04730/2020).

## References

- [1] J. H. O. U. K. V. Johannes Lindvall, “Post-construction production assessment of wind farms - Assessment and optimization of the energy production of operational wind farms: Part 1,” ENERGI FORSK, 2016.
- [2] R. Gelaro, W. McCarty, M.J. Suarez, R. Todling, A. Molod, L. Takacs, C. Randles, A. Darmanov, M.G. Bosilovich, R. Reichle, K. Wargan, L. Coy, R. Cullather, C. Draper, S. Akella, V. Buchar, A. Conaty, A. da Silva, W. Gu, G.-K. Kim, R. Koster, R. Lucchesi, D. Merkova, J.E. Nielsen, G. Partyka, S. Pawson, W. Putman, M. Rienecker, S.D. Schubert, M. Sienkiewicz, B. Zhao, The modern-era retrospective analysis for research and applications, version 2 (MERRA-2), *J. Clim.* 30 (2017) 5419e5454
- [3] R. McKenna, S. Pfenninger, H. Heinrichs, J. Schmidt, I. Staffell, C. Bauer, K. Gruber, A. N. Hahmann, M. Jansen, M. Klingler, N. Landwehr, X. Guo Larsen, J. Lilliestam, B. Pickering, M. Robinius, T. Trondle, O. Turkovska, S. Wehrle, J. Michael Weinand, J. Wohland, “High-resolution large-scale onshore wind energy assessments: A review of potential definitions, methodologies and future research needs,” *Renewable Energy*, ELSEVIER, 2022, pp. 659-684.
- [4] H-Y Kim, “Analysis of variance (ANOVA) comparing means of more than two groups”, Department of Dental Laboratory Science and Engineering, College of Health Science & Department of Public Health Science, Korea, 2014.
- [5] N. M. Razali, B. W. Yap, “Power Comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors and Anderson-Darling Tests”, *Journal of Statistical Modeling and Analytics*, Vol.2 No.1, 21-33, 2011.
- [6] Durbin, J., & Watson, G. S., “Testing for Serial Correlation in Least Squares Regression: I”. *Biometrika*, 409-428, 1950.
- [7] Breusch, T. S. and Pagan A. R., “A simple test for heteroscedasticity and random coefficient variation”, *Econometrica* 47, 1287-1294, 1979.