# Solar Energy Prediction Using Weather Forecasts: A Comparative Evaluation of Machine Learning Models

Nancy W. Riad[1*] and Nour A. Mohamed [1]

## ABSTRACT

*Abstract—* *This study advances solar energy forecasting by developing and evaluating a robust framework that integrates high-resolution weather data with machine learning models. The unpredictable nature of solar energy generation, primarily due to intermittent weather conditions, poses a significant challenge for its efficient integration into power grids. To address this, our research leverages critical meteorological parameters, including Global Horizontal Irradiance (GHI), ambient temperature, relative humidity, cloud cover, and precipitation, to predict photovoltaic (PV) energy output. We perform a comprehensive comparative evaluation of five distinct machine learning algorithms: Linear Regression, Decision Tree, Lasso Regression, Gradient Boosting, and XGBoost, to determine the most effective model for this task. The models were trained on a comprehensive historical dataset comprising 196,776 records, which incorporates both meteorological inputs and corresponding solar energy measurements. The results of our evaluation reveal that XGBoost achieves a superior predictive accuracy of 94%, enabling reliable long-term forecasts and enhanced grid management. The framework demonstrates that leveraging meteorological data significantly boosts prediction performance, thereby supporting sustainable and reliable energy integration. Our findings underscore the efficacy of ensemble learning methods, particularly XGBoost, in capturing the complex, non-linear dependencies inherent in solar energy data.*

Keywords: Solar energy forecasting, photovoltaic output prediction, machine learning, meteorological parameters, global horizontal irradiance (GHI), gradient boosting, XGBoost, renewable energy integration, Long-term forecasting, ensemble learning

## I.    INTRODUCTION

In recent years, the demand for renewable energy sources has grown significantly due to increasing concerns over environmental sustainability and the depletion of fossil fuels. Among various renewable sources, solar energy has emerged as one of the most promising and widely adopted alternatives, primarily due to its abundance and environmental friendliness. [1] The International Energy Agency (IEA) highlighted in its 2023 report the crucial role of solar power in the global energy transition.

However, the efficient integration of solar energy into the power grid presents a significant challenge, primarily due to the intermittent and unpredictable nature of solar irradiance. Variations in weather conditions such as cloud cover, temperature, humidity, and wind speed can significantly impact solar energy generation. Hence, accurate forecasting of solar energy production is essential for optimizing energy management systems, ensuring grid stability, and improving the overall performance of solar power systems. [2]

Without precise forecasts, grid operators face difficulties in balancing supply and demand, which can lead to operational uncertainties and potential grid instability. This unpredictability

---

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

necessitates the development of sophisticated predictive models to support informed decision-making for energy storage, load balancing, and dynamic grid integration.

Traditional statistical methods, while foundational, have shown limitations in capturing the complex, nonlinear patterns inherent in solar energy data. The dynamic and multifaceted relationship between atmospheric conditions and solar power output requires more advanced techniques. In contrast, Artificial Intelligence (AI) algorithms—particularly Machine Learning (ML) and Deep Learning (DL) techniques—have demonstrated superior capabilities in modeling such complex relationships. [3] Recent studies have successfully applied a variety of machine learning models, including Support Vector Machines (SVMs), Random Forests (RF), and Long Short-Term Memory networks (LSTM), to energy prediction tasks with promising results. [4] These methods can analyze vast datasets and identify subtle correlations that are often missed by conventional approaches.

Forecasting photovoltaic (PV) energy output has attracted significant research interest over the past decade due to its critical role in renewable energy integration and grid management. Early studies mainly relied on statistical techniques such as autoregressive models and linear regression, but these approaches struggled to capture the nonlinear behavior of solar energy generation under varying meteorological conditions.

In summary, prior literature demonstrates both the potential and the limitations of existing PV forecasting approaches. While statistical and machine learning methods have significantly advanced the field, gaps remain in fully leveraging detailed meteorological data to improve long-term and multi-horizon prediction accuracy. This study addresses these gaps by developing a robust forecasting framework that integrates high-resolution weather features with advanced machine learning algorithms.

Table 1 provides a comparative summary of the selected literature. The table contrasts the focus, methodologies, and key findings of each work, thereby illustrating how the current study builds upon and extends previous research.

| Study | Method / Model | Input Features | Dataset | Key Findings |
|---|---|---|---|---|
| Marzouk (2025) [1] | Policy/Framework analysis | Clean energy indicators (aggregate rating) | IEA TCEP 2023 Report | Proposed single aggregate rating for tracking clean energy progress. |
| Pedro & Coimbra (2012) [2] | Statistical (no exogenous inputs) | Historical PV output only | Solar plant data (2012) | Highlighted limitations of forecasting without weather data. |
| Voyant et al. (2017) [3] | Machine Learning (ANN, SVM, etc.) | Meteorological + solar irradiance | Multiple datasets (review) | Showed ML models outperform traditional approaches. |
| Ahmed et al. (2020) [4] | Review + Optimization | PV + weather parameters | 124 studies analyzed | Emphasized optimization methods and advanced forecasting techniques. |

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

| | | | | |
|---|---|---|---|---|
| Proposed Study | XGBoost, Gradient Boosting, Decision Trees, Lasso Regression, Linear Regression | High-resolution weather features (GHI, temp, humidity, cloud cover, precipitation, etc.) | 196,776 records (real-world dataset) | Improved accuracy and robustness; addresses gaps in detailed meteorological integration. |

TABLE I: Comparative summary of past and recent studies on PV energy forecasting

This paper investigates the use of AI algorithms in predicting solar energy generation. The main objective is to evaluate and compare the performance of different AI-based models using real-world meteorological and solar radiation data. By identifying the most accurate and robust model, this study aims to contribute to the development of a more intelligent and reliable solar energy forecasting system. The integration of high-resolution weather data, such as Global Horizontal Irradiance (GHI), temperature, humidity, and cloud cover, is a key focus of this research, as it has been shown to significantly boost prediction performance and support the sustainable integration of solar energy.

## II.    METHODOLOGY

This study develops a sophisticated forecasting framework for photovoltaic (PV) energy output by integrating high-resolution weather prediction variables. The precise anticipation of solar generation is paramount for optimizing grid resource allocation, mitigating operational uncertainties, and enhancing the efficiency of renewable energy systems. We leverage critical meteorological parameters, including global horizontal irradiance (GHI), ambient temperature, relative humidity, cloud cover fraction, and precipitation, to engineer data-driven models that capture complex non-linear dependencies between atmospheric dynamics and PV power generation.

Our machine learning architecture enables adaptive forecasting across multiple time horizons. This AI-powered approach significantly improves the operational reliability and economic feasibility of solar-dominated power infrastructures [5,6]

Our machine learning architecture evaluates a selection of models, including XGBoost, Gradient Boosting, Decision Trees, and regression variants. The model was trained on a comprehensive historical dataset of 196,776 records, which integrates both meteorological inputs and corresponding solar energy measurements. [7] This rich dataset allows the model to learn complex patterns and dependencies that influence solar energy generation.

The dataset includes a diverse set of features categorized as follows:

- Temporal features: Timestamp, hour, month.

- Energy measurements: Energy delta (Wh).

- Solar metrics: Global Horizontal Irradiance (GHI), sunlight indicators (is Sun, sunlight time, day length).

- Weather parameters: Temperature, pressure, humidity, wind speed, precipitation (rain_lh, snow_1h), cloud cover (clouds all), weather type.

- Derived ratios: Sunlight time to day length ratio.

---

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

To construct an effective predictive framework, we developed a machine learning model that leverages structured weather forecast data to accurately estimate solar energy generation from photovoltaic (PV) panels. The model was trained on extensive historical datasets comprising both meteorological variables and corresponding energy output measurements, enabling it to capture underlying patterns and complex relationships that drive solar energy production.[8]

The modelling workflow encompasses several critical stages, including data preprocessing—such as addressing missing values and normalizing inputs—feature selection to identify the most impactful weather parameters, followed by model training and rigorous evaluation. The selected algorithm establishes a functional mapping between input features, including temperature, humidity, cloud cover, and solar radiation, and the target output variable, namely the energy generated by the solar panels.

The primary objective is to maximize predictive accuracy to facilitate real-time operational decision-making and long-term strategic planning for solar power installations. Enhanced forecasting precision supports improved management of energy storage systems, demand prediction, and scalability considerations, thereby advancing the efficiency and sustainability of renewable energy infrastructure.[9]

### A. *Linear Regression*

Linear regression is a classical supervised machine learning technique that utilises labelled datasets to learn the optimal linear relationship between input features and the target variable. It aims to find the best-fitting linear function that can accurately predict outcomes on new, unseen data. [10]

This method fundamentally models the dependency of a response variable—here, the solar energy output—on one or more predictor variables, such as weather parameters, by fitting a linear equation to the observed data. The underlying assumption of the model is that the relationship between the dependent and independent variables is linear, which can be mathematically represented as: [11]

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + \varepsilon \qquad (1)$$

Where:

- y = predicted energy output

- $\beta_0$ = intercept

- $\beta_1 \ldots \beta_n$ = coefficients for each feature

- $x_1 \ldots x_{1n}$ = input features (GHI, temperature, etc.)

- $\varepsilon$ = error term.

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt
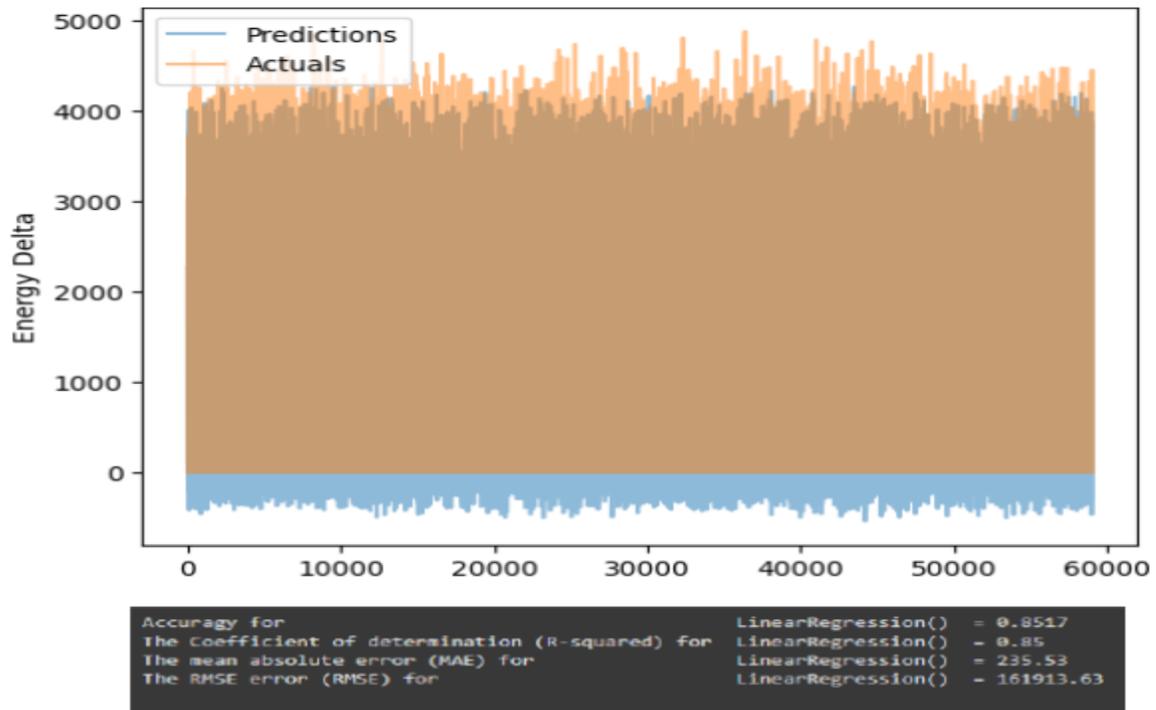
[*] n.wadie@sha.edu.eg

Fig.1: Accuracy of Linear Regression

The algorithm was evaluated on the dataset, achieving an accuracy of 85%. As shown in Fig.1, this performance highlights its effectiveness in modelling the relationship between weather variables and solar energy output, serving as a reliable baseline for further improvement.

## B. Decision tree

Decision trees are supervised learning algorithms that recursively partition data into subsets based on feature values, forming a tree-like structure where each internal node represents a decision rule and each leaf node corresponds to a prediction. This splitting process aims to maximize the homogeneity of the resulting subsets according to a chosen criterion, commonly information gain or impurity reduction. [12]

In the context of solar energy forecasting, decision trees partition weather-related variables—such as solar irradiance and cloud cover—using binary splits that optimize measures like Gini impurity or entropy. For example, a node might split the data based on a condition like:

$$Irradiance > 500 \, W/m2$$

separating days with high solar energy production from those with lower output.

Key parameters influencing the model performance include:

• Max Depth ($dmax$): Limits the maximum depth of the tree to control overfitting, especially important when handling noisy weather data.

• Minimum Samples per Leaf ($nmin$): Specifies the minimum number of samples required

---

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

To form a leaf node, preventing over-segmentation caused by rare or outlier weather events.

The tree chooses splits by maximizing the reduction in impurity, often measured by Gini impurity. The information gain from a split is calculated as:

$$\left(G(left) + {N(right)}/{N}\ G(right)\ {N(left)}/{N}\right) - G(parent) = \Delta I \quad (2)$$

where:

$\Delta I$ represents the information gain from making a particular split in the data.

$G$ (.) is the impurity measure of a node,

$N$ is the total number of samples in the parent node,

$Nleft$, $Nright$ are the numbers of samples in the left and right child nodes, respectively

This process continues recursively until stopping criteria such as max depth or minimum samples per leaf are met, resulting in a model capable of capturing nonlinear relationships in solar energy data with interpretability and efficiency. [13]
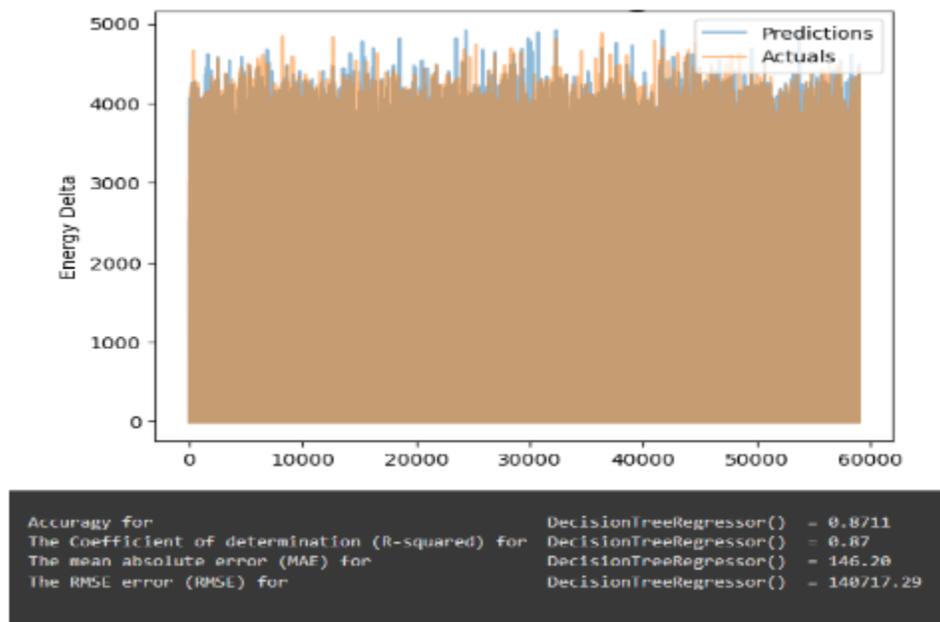


Fig.2: Accuracy of Decision Tree

The algorithm was evaluated on the same dataset, achieving an accuracy of 87%. As shown in Fig.2, this performance highlights its effectiveness in modelling the relationship between weather variables and solar energy output, serving as a reliable baseline for further improvement.

### C. Lasso Regression (L1 Regularization)

The Lasso (Least Absolute Shrinkage and Selection Operator) algorithm is a regularized linear regression technique employed in this study to predict solar energy output based on weather data. Lasso enhances traditional linear regression by incorporating a penalty term that constrains the

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

magnitude of regression coefficients, effectively reducing model complexity and preventing overfitting. Specifically, Lasso minimizes the residual sum of squares while adding an L1 regularization term proportional to the absolute values of the coefficients. This regularization not only improves the generalization ability of the model but also performs feature selection by shrinking less important coefficients towards zero. [14] Mathematically, the Lasso objective function can be expressed as:

$$min_\beta\{\sum_{i=1}^{n}(yi - \beta_0 - \sum_{j=1}^{p}\beta_j\,x_{ij})^2 + \lambda\sum_{j=1}^{p}|\,\beta j\,|\}\qquad(3)$$

where:

$y_i$ = Solar energy output,

$xij$ = Weather features (e.g., irradiance, temperature),

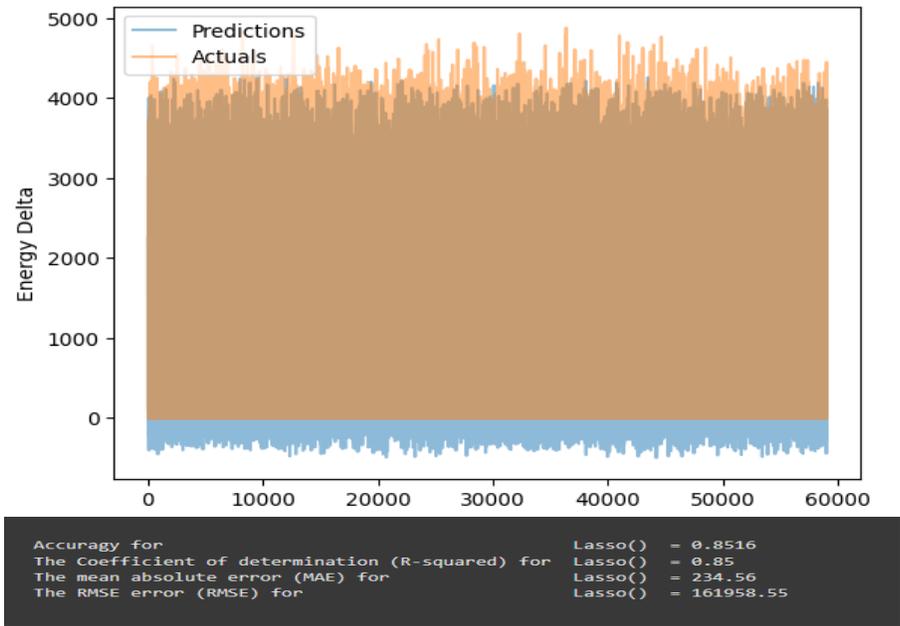$\lambda$ = Penalty strength (tuned via cross-validation).



Fig.3: Accuracy of Lasso Regression

The algorithm was evaluated on the same dataset, achieving an accuracy of 85%. As shown in Fig.3, this performance highlights its effectiveness in modelling the relationship between weather variables and solar energy output, serving as a reliable baseline for further improvement.

### D. Gradient Boosting

Gradient Boosting is an ensemble learning technique that builds predictive models in a sequential manner, where each new model aims to correct the errors of its predecessor by minimizing a specified loss function, such as mean squared error (MSE) or cross-entropy. At each iteration, the algorithm calculates the gradient of the loss function concerning the current model's predictions and fits a new weak learner, typically a shallow decision tree, to approximate this gradient. The updated ensemble prediction at iteration m is given by

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt
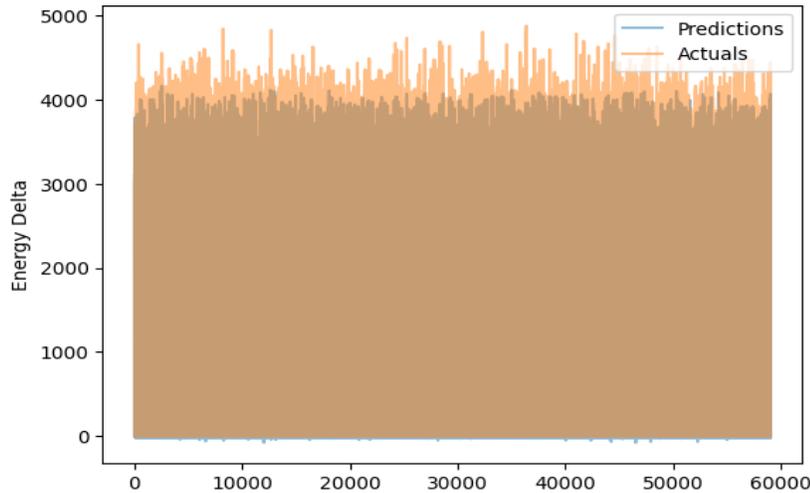
[*] n.wadie@sha.edu.eg

$$F_{m+1}(x) = F_m(x) + \gamma h_m(x)F_m + 1(x) = F_m(x) + \gamma h_m(x) \quad (4)$$

Where:

$h_m(x)$ corrects errors from prior model

$F_m$. Optimized via gradient descent on loss functions like MAE (robust to solar outliers).[15]



```
Accuracy for                                        GradientBoostingRegressor()  = 0.9159
The Coefficient of determination (R-squared) for   GradientBoostingRegressor()  = 0.92
The mean absolute error (MAE) for                   GradientBoostingRegressor()  = 134.46
The RMSE error (RMSE) for                           GradientBoostingRegressor()  = 91837.82
```

Fig.4: Accuracy of Gradient Boosting

The algorithm was evaluated on the same dataset, achieving an accuracy of 92%. As shown in Fig.4, this performance highlights its effectiveness in modelling the relationship between weather variables and solar energy output, serving as a reliable baseline for further improvement.

### E.  XGBoost algorithm

XGBoost uses decision trees as its base learners and combines them sequentially to improve the model's performance. Each new tree is trained to correct the errors made by the previous tree, a process known as boosting. It features built-in parallel processing, enabling the rapid training of models on large datasets.

• Start with a base learner: The first model decision tree is trained on the data. In regression tasks, this base model simply predicts the average of the target variable.

• Calculate the errors: After training the first tree, the errors between the predicted and actual values are calculated.

• Train the next tree: The next tree is trained on the errors of the previous tree. This step attempts to correct the errors made by the first tree.

• Repeat the process: This process continues with each new tree, trying to correct the errors of the previous trees until a stopping criterion is met.

• Combine the predictions: The final prediction is the sum of the predictions from all the trees. [16]

---

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

Mathematically, the model can be represented as:

$$y^{\wedge}_i = \sum_{k=1}^{k} f_k(x_i) \quad (5)$$

Where:

$f_k$ represents the kth decision tree

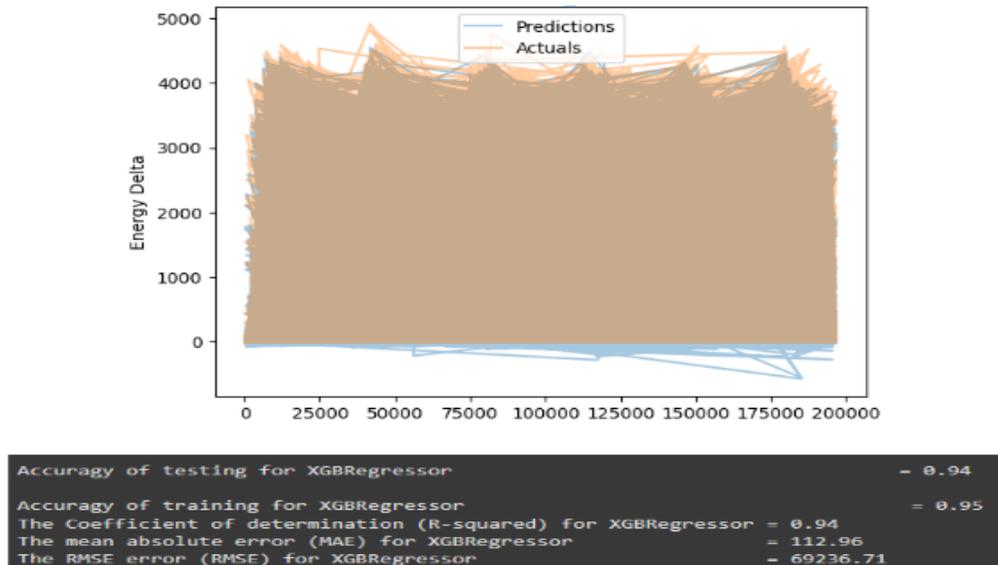$y^{\wedge}\,i$ is the final predicted value, for instance $i$



Fig.5: Accuracy of XGBoosting

The algorithm was evaluated on the same dataset, achieving an accuracy of 94%. As shown in Fig.5, this performance highlights its effectiveness in modelling the relationship between weather variables and solar energy output, serving as a reliable baseline for further improvement.

## III.    RESULTS & DISCUSSION

The results of the comparative evaluation of the five machine learning algorithms are summarized in the table below. The analysis is based on model accuracy, Mean Absolute Error (MAE), and Root Mean Square Error (RMSE).

| Algorithm | Accuracy (%) | MAE | RMSE |
|---|---|---|---|
| Linear Regression | 85.17 | 235.53 | 161913.63 |
| Decision Tree | 87.11 | 146.20 | 140717.29 |
| Lasso Regression | 85.16 | 234.56 | 161958.55 |
| Gradient Boosting | 91.59 | 134.46 | 91837.82 |

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

| | | | |
|---|---|---|---|
| XGBoosting | 94.00 | 112.96 | 69236.71 |

Table II.: A Comparative Study of Algorithms in Photovoltaic Power Prediction

The results clearly indicate that the ensemble learning methods, XGBoost and Gradient Boosting, significantly outperformed the single-model approaches (Linear Regression, Lasso, and Decision Tree). XGBoost achieved the highest accuracy of 94%, which confirms its superior ability to capture complex, non-linear relationships in energy usage data. Its performance, coupled with a lower Mean Absolute Error (MAE) and Root Mean Square Error (RMSE), indicates a more robust and reliable forecasting capability. The Gradient Boosting model also performed exceptionally well, with an accuracy of 92%, demonstrating the general strength of boosting techniques for this task.

In contrast, the Decision Tree model showed a moderate accuracy of 87%, suggesting that a single tree's capacity is limited without the boosting framework. The linear models, Linear Regression and Lasso Regression, both achieved the lowest accuracy of 85%. This supports the conclusion that simpler linear models are insufficient for accurately modeling the high variability and complex patterns in solar energy consumption. The difference in performance highlights that the complex dependencies between meteorological variables and energy output are better captured by more sophisticated, non-linear models.

## V. CONCLUSION

This study aimed to evaluate the performance of various machine learning algorithms for energy consumption prediction. The comparison, based on model accuracy, reveals that XGBoost achieved the highest accuracy of 94%, making it the most effective model for capturing complex, non-linear relationships in energy usage data. Gradient Boosting followed closely with 92%, also demonstrating strong predictive power due to its ensemble nature. The Decision Tree model showed moderate accuracy at 87%, suggesting its limited capacity without the boosting framework.

Both Lasso Regression and Linear Regression achieved the lowest accuracy of 85%, indicating that simpler linear models may not be sufficient for accurately modelling the variability in energy consumption patterns.

In conclusion, ensemble learning methods, particularly XGBoost, provide superior performance in energy prediction tasks and are better suited for applications that require high accuracy and robustness in forecasting energy demand.

### ACKNOWLEDGMENT

### REFERENCES

[1] O.A. Marzouk, Summary of the 2023 report of TCEP (tracking clean energy progress) by the International Energy Agency (IEA), and proposed process for computing a single aggregate rating, E3S Web Conf. 601 (2025) 1–6. https://doi.org/10.1051/e3sconf/202560100048

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg

[2] H.T.C. Pedro, C.F.M. Coimbra, Assessment of forecasting techniques for solar power production with no exogenous inputs, Solar. Energy 86 (2012) 2017–2028. https://doi:10.1016/j.solener.2012.04.004

[3] C. Voyant, G. Notton, S. Kalogirou, M.-L. Nivet, C. Paoli, F. Motte, A. Fouilloy, Machine learning methods for solar radiation forecasting: A review, Renew. Energy 105 (2017) 569–582. https://doi:10.1016/j.renene.2016.12.095

[4] R. Ahmed, V. Sreeram, Y. Mishra, M.D. Arif, A review and evaluation of the state-of-the-art in PV solar power forecasting: Techniques and optimization, Renew. Sustain. Energy Rev. 124 (2020) 109792. https://doi.org/10.1016/j.rser.2020.109792

[5] W.-C. Tsai, C.-H. Chen, Y.-H. Chang, C.-M. Lin, A review of state-of-the-art and short-term forecasting models for solar PV power generation, *Energies* 16 (2023) 5436. https://doi.org/10.3390/en16145436

[6] F. Wang, Z. Liu, X. Zhang, Y. Wang, K. Zhao, Generative adversarial networks and convolutional neural networks-based weather classification model for day-ahead short-term photovoltaic power forecasting, *Energy Convers. Manag.* 181 (2019) 443–462. https://doi:10.1016/j.enconman.2018.11.074

[7] S. Emami, Renewable Energy and Weather Conditions [dataset], Kaggle, 2023. https://www.kaggle.com/datasets/samanemami/renewable-energy-and-weather-conditions.

[8] A. Yona, T. Senjyu, T. Funabashi, H. Sekine, Application of neural network to 24-hour-ahead generating power forecasting for PV system, *Proc. IEEE Power Energy Soc. Gen. Meet.* (2008) 1–6. https://doi:10.1541/ieejpes.128.33

[9] A.K. Ozcanli, F. Yaprakdal, M. Baysal, Deep learning methods and applications for electrical power systems: A comprehensive review, *Int. J. Energy Res.* 44 (2020) 7136–7157. https://doi:10.1002/er.5331

[10] GeeksforGeeks Machine Learning–*Linear Regression.* https://www.geeksforgeeks.org/machine-learning/ml-linear-regression, 2023.

[11] J. Antonanzas, N. Osorio, R. Escobar, R. Urraca, F.J. Martinez-de-Pison, F. Antonanzas-Torres, Review of photovoltaic power forecasting, *Sol. Energy* 136 (2016) 78–111. https://doi:10.1002/er.5331

[12] I. D. Mienye, N. Jere, A survey of decision trees: concepts, algorithms, and applications, *IEEE Access* 12 (2024) 86716–86727.https://doi:10.1109/ACCESS.2024.3416838

[13] L.Wilkinson, Classification and regression trees, in: *Systat* 11, 2004, pp. 35–56.

[14] W. Chen, et al., The Prediction Performance Analysis of the Lasso Model with Convex Non-Convex Sparse Regularization, Algorithms 18 (2025) 195. https://doi.org/10.3390/a18040195

[15] E. A. Tuncar, Ş. Sağlam, B. Oral, A review of short-term wind power generation forecasting methods in recent technological trends, Energy Rep. 12 (2024) 197–209. https://doi:10.1016/j.egyr.2024.06.006

[16] A. Groß, A. König, A. Waczowicz, M. S. Hossain, A. Monti, Comparison of short-term electrical load forecasting methods for different building types, *Energy Informatics* 4 (Suppl 3) (2021) 13. https:// doi:10.1186/s42162-021-00172-6

[1] Department of Communications and Computer Engineering, Higher Institute of Engineering (HIE), El-Shorouk Academy, El-Shorouk City 11837, Egypt

[*] n.wadie@sha.edu.eg