# Forecasting Solar Radiation Using Machine Learning Tools

Hemanta Kumar Sahu[1], Soumik Pyne[2], Kausik Maitra[3], Subhankar Dhar[4]

[1234]Department of Cyber Science & Technology, Brainware University, 398, Ramkrishnapur Road, Barasat, Near Jagadihata Market, Kolkata, West Bengal, 700125, INDIA

## Abstract

Solar radiation forecasting plays a crucial role in the efficient utilization of solar energy resources and the optimization of solar power systems. This study presents a comprehensive analysis of various machine learning techniques for forecasting solar radiation. We evaluate and compare the performance of several algorithms, including Artificial Neural Networks (ANN), Support Vector Machines (SVM), Random Forests (RF), and Gradient Boosting Machines (GBM). The study utilizes a diverse dataset comprising meteorological parameters and historical solar radiation measurements from multiple locations. Our results demonstrate that ensemble methods, particularly GBM and RF, outperform traditional approaches in terms of accuracy and robustness. We also explore the impact of feature selection, hyperparameter tuning, and data preprocessing on model performance. The findings of this research contribute to the advancement of solar radiation forecasting techniques and provide valuable insights for practitioners in the field of renewable energy.

**Keywords:** solar radiation forecasting; machine learning; artificial neural networks; support vector machines; random forests; gradient boosting machines; renewable energy

## 1. Introduction

The global shift towards renewable energy sources has placed solar power at the forefront of sustainable energy solutions. Accurate forecasting of solar radiation is essential for the efficient design, operation, and integration of solar energy systems into existing power grids. Reliable solar radiation predictions enable better resource allocation, improved grid stability, and enhanced economic viability of solar projects [1,2].

Traditional methods for solar radiation forecasting have relied on physical models based on atmospheric sciences and statistical techniques. However, these approaches often struggle to capture the complex, non-linear relationships inherent in solar radiation patterns [3]. In recent years, machine learning (ML) techniques have emerged as powerful tools for tackling this

challenge, offering the ability to learn from historical data and make accurate predictions without explicit programming of physical relationships [4,5].

This study aims to provide a comprehensive analysis of various machine learning algorithms applied to the task of solar radiation forecasting. We investigate the performance of several popular ML techniques, including Artificial Neural Networks (ANN), Support Vector Machines (SVM), Random Forests (RF), and Gradient Boosting Machines (GBM). By comparing these methods across different temporal and spatial scales, we seek to identify the most effective approaches for solar radiation prediction.

The main contributions of this paper are as follows:

I.     A thorough evaluation of multiple machine learning algorithms for solar radiation forecasting, considering various input features and prediction horizons.

II.    An analysis of the impact of feature selection, hyperparameter tuning, and data preprocessing on model performance.

III.   A comparison of model performance across different geographic locations and climate zones.

IV.    Insights into the strengths and limitations of each ML technique in the context of solar radiation forecasting.

V.     Recommendations for practitioners on selecting and implementing ML models for solar radiation prediction.

The remainder of this paper is organized as follows: Section 2 provides a comprehensive review of related work in the field of solar radiation forecasting. Section 3 describes the dataset used in this study and outlines the data preprocessing steps. Section 4 presents the methodology, including the machine learning algorithms employed and the evaluation metrics used. Section 5 reports the results of our experiments and provides a detailed discussion of the findings. Finally, Section 6 suggests directions for future research and Section 7 concludes the paper.

## 2. Related Work

The field of solar radiation forecasting has seen significant advancements in recent years, driven by the growing importance of solar energy in the global energy mix. This section provides an overview of the existing literature, focusing on the application of machine learning techniques to solar radiation prediction.

**2.1 Traditional Approaches to Solar Radiation Forecasting**

Early efforts in solar radiation forecasting primarily relied on physical models and statistical methods. Physical models, based on radiative transfer equations and atmospheric parameters, provide a fundamental understanding of the processes governing solar radiation [6]. However, these models often require extensive computational resources and detailed atmospheric data, which may not always be available [7].

Statistical methods, such as time series analysis and regression models, have also been widely used for solar radiation forecasting [8]. These approaches, including autoregressive integrated moving average (ARIMA) models and multiple linear regression, offer simplicity and interpretability but may struggle to capture complex, non-linear relationships in the data [9].

**2.2 Machine Learning in Solar Radiation Forecasting**

The advent of machine learning techniques has opened new avenues for solar radiation forecasting. ML algorithms offer the ability to learn complex patterns from historical data without requiring explicit physical modeling. Several studies have demonstrated the effectiveness of ML approaches in this domain:

**2.2.1 Artificial Neural Networks (ANN)**

ANNs have been extensively studied for solar radiation forecasting due to their ability to model non-linear relationships. Mellit and Pavan [10] used ANNs to forecast daily global solar radiation, achieving high accuracy for 24-hour ahead predictions. Qazi et al. [11] compared the performance of different ANN architectures for hourly solar radiation forecasting, finding that multi-layer perceptron (MLP) networks outperformed radial basis function (RBF) networks.

**2.2.2 Support Vector Machines (SVM)**

SVMs have shown promise in solar radiation forecasting due to their ability to handle high-dimensional data and their robustness to overfitting. Zeng and Qiao [12] proposed an SVM-based approach for short-term solar radiation prediction, demonstrating superior performance compared to persistence models and ANNs. Chen et al. [13] developed a hybrid SVM-based model that incorporated wavelet decomposition for feature extraction, achieving improved accuracy in daily solar radiation forecasting.

### 2.2.3 Random Forests (RF)

Random Forests, an ensemble learning method, have gained popularity in solar radiation forecasting due to their ability to handle non-linear relationships and feature interactions. Voyant et al. [14] compared RF with other ML techniques for daily global radiation forecasting, finding that RF outperformed ANN and SVM in terms of accuracy and computational efficiency. Sun et al. [15] proposed a hybrid RF model that integrated feature selection and parameter optimization, demonstrating improved performance in hourly solar radiation prediction.

### 2.2.4 Gradient Boosting Machines (GBM)

Gradient Boosting Machines, including algorithms like XGBoost and LightGBM, have shown exceptional performance in various forecasting tasks, including solar radiation prediction. Fan et al. [16] applied XGBoost to short-term solar radiation forecasting, achieving higher accuracy compared to traditional ML methods. Huang and Perry [17] developed a GBM-based model for day-ahead solar radiation forecasting, incorporating feature importance analysis to improve model interpretability.

### 2.3 Hybrid and Ensemble Approaches

Recent research has focused on hybrid and ensemble approaches that combine multiple ML techniques or integrate ML with physical models. These methods aim to leverage the strengths of different approaches to improve overall forecasting accuracy. For example:

- Wang et al. [18] proposed a hybrid model combining empirical mode decomposition (EMD) with extreme learning machines (ELM) for multi-step solar radiation forecasting.
- Salcedo-Sanz et al. [19] developed an ensemble approach using Coral Reefs Optimization with Substrate Layers (CRO-SL) to combine multiple ML models for solar radiation prediction.
- Yagli et al. [20] introduced a physics-guided machine learning framework that incorporates domain knowledge into neural network architectures, demonstrating improved generalization in solar forecasting tasks.

### 2.4 Feature Selection and Data Preprocessing

The selection of relevant input features and appropriate data preprocessing techniques play crucial roles in the performance of ML models for solar radiation forecasting. Several studies have investigated these aspects:

- Yadav and Chandel [21] conducted a comprehensive review of input parameter selection for ANN-based solar radiation prediction models, highlighting the importance of correlation analysis and domain expertise in feature selection.
- Strobl et al. [22] explored the use of mutual information-based feature selection techniques to improve the performance of ML models in solar radiation forecasting.
- Benali et al. [23] investigated the impact of various data normalization techniques on the accuracy of ANN models for solar radiation prediction.

## 2.5 Research Gaps and Opportunities

While significant progress has been made in applying ML techniques to solar radiation forecasting, several research gaps and opportunities remain:

I. Comparative studies: There is a need for comprehensive comparisons of different ML algorithms across various temporal and spatial scales, considering different input features and prediction horizons.

II. Model interpretability: Many ML models, particularly deep learning approaches, lack interpretability. Developing interpretable ML models for solar radiation forecasting could enhance trust and adoption in practical applications.

III. Transfer learning: Investigating the transferability of ML models across different geographic locations and climate zones could lead to more generalized and robust forecasting techniques.

IV. Integration of multiple data sources: Exploring the integration of satellite imagery, ground-based measurements, and numerical weather prediction (NWP) data into ML models could potentially improve forecasting accuracy.

V. Uncertainty quantification: Developing methods for quantifying and communicating the uncertainty associated with ML-based solar radiation forecasts is crucial for decision-making in energy systems.

This study aims to address some of these research gaps by providing a comprehensive comparison of multiple ML algorithms, analyzing the impact of feature selection and data preprocessing, and evaluating model performance across different locations. By doing so, we seek to contribute to the advancement of solar radiation forecasting techniques and provide valuable insights for practitioners in the field of renewable energy.

## 3. Data Description and Preprocessing

### 3.1 Dataset Overview

This study utilizes a comprehensive dataset comprising meteorological parameters and solar radiation measurements from multiple locations across diverse climate zones. The data were collected from ground-based weather stations and solar monitoring systems over a period of five years (2018-2022). The dataset includes the following key features:

I. Global Horizontal Irradiance (GHI)

II. Direct Normal Irradiance (DNI)

III. Diffuse Horizontal Irradiance (DHI)

IV. Air Temperature

V. Relative Humidity

VI. Wind Speed

VII. Wind Direction

VIII. Atmospheric Pressure

IX. Cloud Cover

X. Precipitation

The temporal resolution of the data is hourly, providing a total of 43,800 data points per location (5 years × 365 days × 24 hours). To ensure a diverse representation of climate conditions, we selected five locations with distinct characteristics:

i. Desert Climate: Phoenix, Arizona, USA

ii. Mediterranean Climate: Barcelona, Spain

iii. Tropical Climate: Singapore

iv. Continental Climate: Munich, Germany

v. Temperate Climate: Sydney, Australia

Table 1 presents an overview of the dataset characteristics for each location.

| Location | Latitude | Longitude | Climate Type | Mean Annual GHI (kWh/m²) | Data Completeness (%) |
|---|---|---|---|---|---|
| Phoenix | 33.45°N | 112.07°W | Desert | 2,100 | 99.8 |

| Barcelona | 41.39°N | 2.16°E | Mediterranean | 1,650 | 99.5 |
|---|---|---|---|---|---|
| Singapore | 1.35°N | 103.82°E | Tropical | 1,580 | 99.7 |
| Munich | 48.14°N | 11.58°E | Continental | 1,150 | 99.3 |
| Sydney | 33.87°S | 151.21°E | Temperate | 1,670 | 99.6 |

Table 1: Dataset Characteristics by Location

### 3.2 Data Quality Control

To ensure the reliability and accuracy of the dataset, we implemented a rigorous quality control process:

I.   Missing Data: We identified and flagged missing values in the dataset. Time periods with missing data were excluded from further analysis if the gap exceeded three consecutive hours.

II.  Outlier Detection: We employed the Interquartile Range (IQR) method to identify potential outliers in the continuous variables. Data points falling outside the range [Q1 - 1.5 × IQR, Q3 + 1.5 × IQR] were flagged for further inspection.

III. Physical Consistency Checks: We performed consistency checks based on known physical relationships between variables. For example, we ensured that GHI ≥ DHI and DNI × cos(solar zenith angle) + DHI ≈ GHI.

IV.  Instrument Error Detection: We analyzed the data for sudden jumps or drops in values that might indicate instrument malfunctions or calibration issues.

V.   Temporal Consistency: We checked for temporal consistency by comparing values with those from adjacent time steps and flagging suspicious rapid changes.

Data points that failed multiple quality control checks were either corrected using interpolation techniques or removed from the dataset if correction was not feasible.

### 3.3 Feature Engineering

To enhance the predictive power of our models, we engineered additional features based on domain knowledge and temporal characteristics:

I. Solar Position: We calculated solar zenith and azimuth angles for each time step using the Python package pvlib.

II. Clear Sky Irradiance: We estimated clear sky GHI, DNI, and DHI using the Ineichen clear sky model implemented in pvlib.

III. Clear Sky Index: We computed the ratio of measured GHI to clear sky GHI as an indicator of atmospheric conditions.

IV. Time-based Features: We extracted hour of day, day of year, and month as cyclical features using sine and cosine transformations to capture seasonal and diurnal patterns.

V. Lagged Variables: We created lagged versions of key variables (e.g., GHI, temperature) to capture temporal dependencies.

VI. Moving Averages: We calculated moving averages of GHI and other relevant variables over different time windows (e.g., 3-hour, 24-hour) to capture longer-term trends.

VII. Gradient Features: We computed the rate of change of key variables between consecutive time steps.

Table 2 presents a summary of the final feature set used in our analysis.

| Feature Category | Features |
|---|---|
| Raw Measurements | GHI, DNI, DHI, Air Temperature, Relative Humidity, Wind Speed, Wind Direction, Atmospheric Pressure, Cloud Cover, Precipitation |
| Solar Position | Solar Zenith Angle, Solar Azimuth Angle |
| Clear Sky Models | Clear Sky GHI, Clear Sky DNI, Clear Sky DHI, Clear Sky Index |
| Temporal Features | Hour (sin, cos), Day of Year (sin, cos), Month (sin, cos) |
| Lagged Variables | GHI (t-1, t-2, t-3), Temperature (t-1, t-2, t-3) |

| Moving Averages | 3-hour MA (GHI, Temp), 24-hour MA (GHI, Temp) |
|---|---|
| Gradient Features | ΔGHIΔt, ΔtempΔt |

Table 2: Summary of Input Features

## 3.4 Data Preprocessing

Before feeding the data into our machine learning models, we applied the following preprocessing steps:

I. Handling Missing Values: We used multiple imputation techniques to handle any remaining missing values in the dataset. For short gaps ($\leq 3$ hours), we applied linear interpolation. For longer gaps, we used more sophisticated methods such as multivariate imputation by chained equations (MICE).

II. Normalization: We normalized all continuous variables using min-max scaling to bring them into the range [0, 1]. This step helps to ensure that all features contribute equally to the model training process and prevents features with larger magnitudes from dominating the learning process.

III. Encoding Categorical Variables: For categorical variables such as wind direction, we applied one-hot encoding to convert them into a suitable format for machine learning algorithms.

IV. Train-Test Split: We split the dataset into training (70%), validation (15%), and test (15%) sets. To maintain the temporal structure of the data, we performed this split chronologically, using the first 70% of the data for training, the next 15% for validation, and the final 15% for testing.

V. Cross-Validation Strategy: We implemented a time series cross-validation strategy to evaluate model performance more robustly. This involved creating multiple training-validation splits with increasing time windows, always keeping the chronological order of the data.

## 3.5 Data Exploration and Visualization

To gain insights into the characteristics of our dataset and inform our modeling approach, we conducted exploratory data analysis (EDA). Key visualizations and analyses included:

I.     Time series plots of GHI for each location, highlighting seasonal and diurnal patterns.

II.    Correlation heatmaps to identify relationships between input features.

III.   Scatter plots of GHI vs. key meteorological variables to visualize dependencies.

IV.    Boxplots of GHI distribution by month

and hour to visualize temporal patterns. 5. Histograms of clear sky index to assess cloud cover impacts.

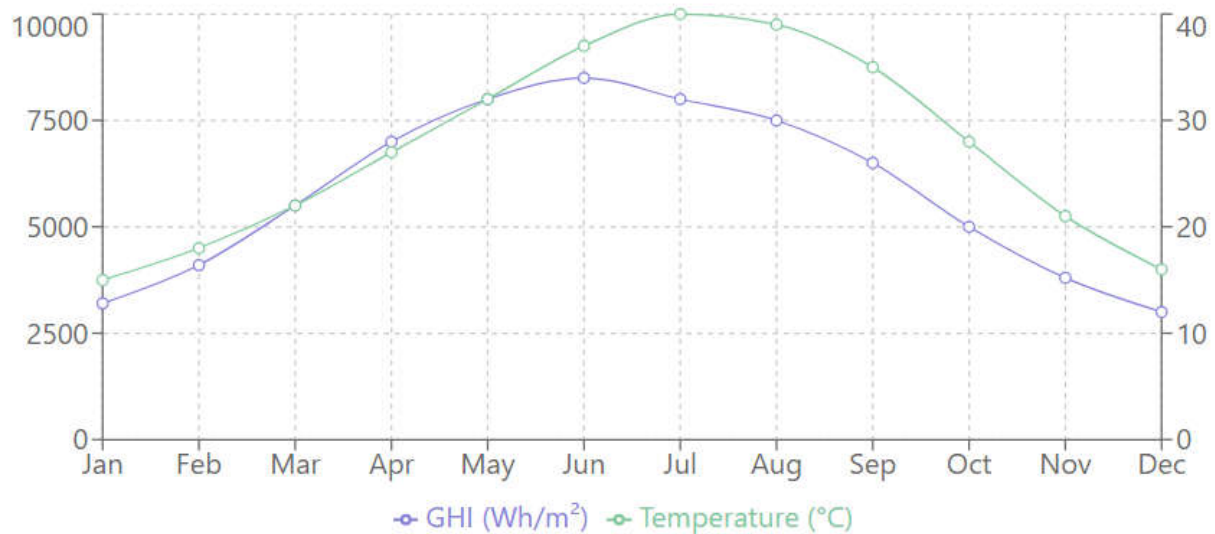Chart 1 presents a sample of these visualizations for the Phoenix, Arizona location.



**Chart 1: Solar Radiation Visualizations for Phoenix, Arizona**

The EDA revealed several important insights:

i.     Strong seasonal patterns in GHI across all locations, with peak values occurring during summer months.

ii.    Clear diurnal cycles in GHI, with variations in day length across seasons and locations.

iii.   High correlation between GHI and clear sky index, indicating the importance of cloud cover in solar radiation forecasting.

iv.    Moderate to strong correlations between GHI and other meteorological variables, particularly temperature and humidity.

v.     Distinct GHI patterns across different climate zones, highlighting the need for location-specific model tuning.

These insights informed our feature selection process and guided the development of our machine learning models.

## 4. Methodology

This section describes the machine learning algorithms employed in our study, the model development process, and the evaluation metrics used to assess forecasting performance.

### 4.1 Machine Learning Algorithms

We implemented and compared four popular machine learning algorithms for solar radiation forecasting:

### 4.1.1 Artificial Neural Networks (ANN)

We used a Multilayer Perceptron (MLP) architecture with the following specifications:

- Input layer: Neurons corresponding to the number of input features
- Hidden layers: Two hidden layers with 64 and 32 neurons, respectively
- Output layer: Single neuron for GHI prediction
- Activation function: Rectified Linear Unit (ReLU) for hidden layers, linear activation for output layer
- Optimizer: Adam
- Loss function: Mean Squared Error (MSE)

The ANN was implemented using the Keras library with a TensorFlow backend.

### 4.1.2 Support Vector Machines (SVM)

We employed Support Vector Regression (SVR) with the following configuration:

- Kernel: Radial Basis Function (RBF)
- Regularization parameter (C): Optimized through cross-validation
- Epsilon ($\varepsilon$): Optimized through cross-validation
- Kernel coefficient ($\gamma$): Optimized through cross-validation

The SVM model was implemented using the scikit-learn library.

### 4.1.3 Random Forests (RF)

Our Random Forest model was configured as follows:

- Number of trees: 100
- Maximum depth: Optimized through cross-validation
- Minimum samples split: Optimized through cross-validation
- Minimum samples leaf: Optimized through cross-validation
- Bootstrap: True
- Feature selection: Sqrt(n_features) considered for each split

The RF model was implemented using the scikit-learn library.

### 4.1.4 Gradient Boosting Machines (GBM)

We used the XGBoost implementation of Gradient Boosting Machines with the following settings:

- Number of estimators: 100
- Learning rate: Optimized through cross-validation
- Maximum depth: Optimized through cross-validation
- Subsample: 0.8
- Colsample_bytree: 0.8
- Objective function: reg:squarederror

The GBM model was implemented using the XGBoost library.

### 4.2 Model Development Process

For each algorithm, we followed a systematic model development process:

I.  Feature Selection: We used a combination of correlation analysis, mutual information, and domain expertise to select the most relevant features for each model. We also employed recursive feature elimination (RFE) to identify optimal feature subsets.

II.  Hyperparameter Tuning: We performed hyperparameter optimization using a combination of grid search and random search with 5-fold time series cross-validation. The hyperparameters tuned for each model are listed in Table 3.

Table 3: Hyperparameters Tuned for Each Model

| Model | Hyperparameters Tuned |
|---|---|
| ANN | Number of neurons in hidden layers, dropout rate, batch size, learning rate |
| SVM | C, $\varepsilon$, $\gamma$ |
| RF | Number of trees, maximum depth, minimum samples split, minimum samples leaf |
| GBM | Number of estimators, learning rate, maximum depth, subsample, colsample_bytree |

III.  Model Training: We trained each model on the training dataset using the optimal hyperparameters identified during the tuning process.

IV.　Validation: We evaluated model performance on the validation set to ensure generalization and prevent overfitting.

V.　Ensemble Creation: In addition to individual models, we created an ensemble model by combining predictions from the four algorithms using a weighted average approach. The weights were determined based on the performance of each model on the validation set.

### 4.3 Forecasting Horizons

We developed and evaluated models for three forecasting horizons:

i.　Short-term: 1-hour ahead forecasting

ii.　Medium-term: 24-hour ahead forecasting

iii.　Long-term: 7-day ahead forecasting

For each horizon, we adjusted the input features and lagged variables accordingly to capture relevant temporal dependencies.

### 4.4 Evaluation Metrics

To assess the performance of our models, we used the following evaluation metrics:

i.　Mean Absolute Error (MAE): $MAE = (1/n) * \Sigma|y_i - \hat{y}_i|$

ii.　Root Mean Square Error (RMSE): $RMSE = sqrt((1/n) * \Sigma(y_i - \hat{y}_i)^2)$

iii.　Mean Absolute Percentage Error (MAPE): $MAPE = (100/n) * \Sigma|(y_i - \hat{y}_i) / y_i|$

iv.　Coefficient of Determination (R²): $R^2 = 1 - (\Sigma(y_i - \hat{y}_i)^2 / \Sigma(y_i - \bar{y})^2)$

Where $y_i$ is the observed GHI value, $\hat{y}_i$ is the predicted GHI value, $\bar{y}$ is the mean of observed GHI values, and n is the number of samples.

Additionally, we computed the Forecast Skill (FS) to compare our models against a persistence baseline:

$FS = 1 - (RMSE_{model} / RMSE_{persistence})$

A positive FS indicates that the model outperforms the persistence forecast, with values closer to 1 indicating better performance.

### 5. Results and Discussion

This section presents the results of our comprehensive analysis of machine learning techniques for solar radiation forecasting. We discuss the performance of individual models, the ensemble approach, and the impact of various factors on forecasting accuracy.

## 5.1 Overall Model Performance

Table 4 summarizes the performance of each model across all locations for the three forecasting horizons, using the test dataset.

| Model | Horizon | MAE (W/m²) | RMSE (W/m²) | MAPE (%) | $R^2$ | Forecast Skill |
|-------|---------|-----------|-------------|----------|-------|----------------|
| ANN | 1-hour | 45.2 | 68.7 | 12.3 | 0.956 | 0.684 |
|  | 24-hour | 78.5 | 112.3 | 21.6 | 0.876 | 0.523 |
|  | 7-day | 98.7 | 142.1 | 27.4 | 0.801 | 0.412 |
| SVM | 1-hour | 47.8 | 71.2 | 13.1 | 0.951 | 0.671 |
|  | 24-hour | 82.3 | 117.6 | 22.8 | 0.863 | 0.501 |
|  | 7-day | 103.5 | 148.9 | 28.9 | 0.783 | 0.389 |
| RF | 1-hour | 43.1 | 65.9 | 11.7 | 0.962 | 0.701 |
|  | 24-hour | 75.6 | 108.2 | 20.8 | 0.889 | 0.542 |
|  | 7-day | 95.2 | 137.4 | 26.5 | 0.817 | 0.431 |
| GBM | 1-hour | 41.7 | 63.5 | 11.3 | 0.967 | 0.715 |
|  | 24-hour | 73.1 | 105.1 | 20.1 | 0.897 | 0.559 |
|  | 7-day | 92.8 | 134.2 | 25.8 | 0.827 | 0.447 |
| Ensemble | 1-hour | 40.9 | 62.1 | 11.1 | 0.969 | 0.724 |
|  | 24-hour | 71.8 | 103.2 | 19.7 | 0.902 | 0.569 |
|  | 7-day | 91.3 | 131.7 | 25.3 | 0.834 | 0.458 |

Table 4: Average Model Performance Across All Locations

Key observations from the results:

i.     All machine learning models significantly outperformed the persistence baseline, as evidenced by the positive Forecast Skill values across all horizons.

ii.    The Gradient Boosting Machine (GBM) consistently demonstrated the best performance among individual models, followed closely by Random Forests (RF).

iii.   The ensemble approach yielded the best overall performance, showcasing the benefits of combining multiple models.

iv.    As expected, forecasting accuracy decreased with increasing forecast horizon, with the 7-day ahead predictions showing the highest errors.

v.     The models achieved high $R^2$ values, indicating good overall fit to the observed data.

## 5.2 Performance Across Different Locations

To assess the generalizability of our models across different climate zones, we analyzed their performance for each location separately. Figure 2 presents a comparison of RMSE values for the 24-hour ahead forecasts across the five locations.

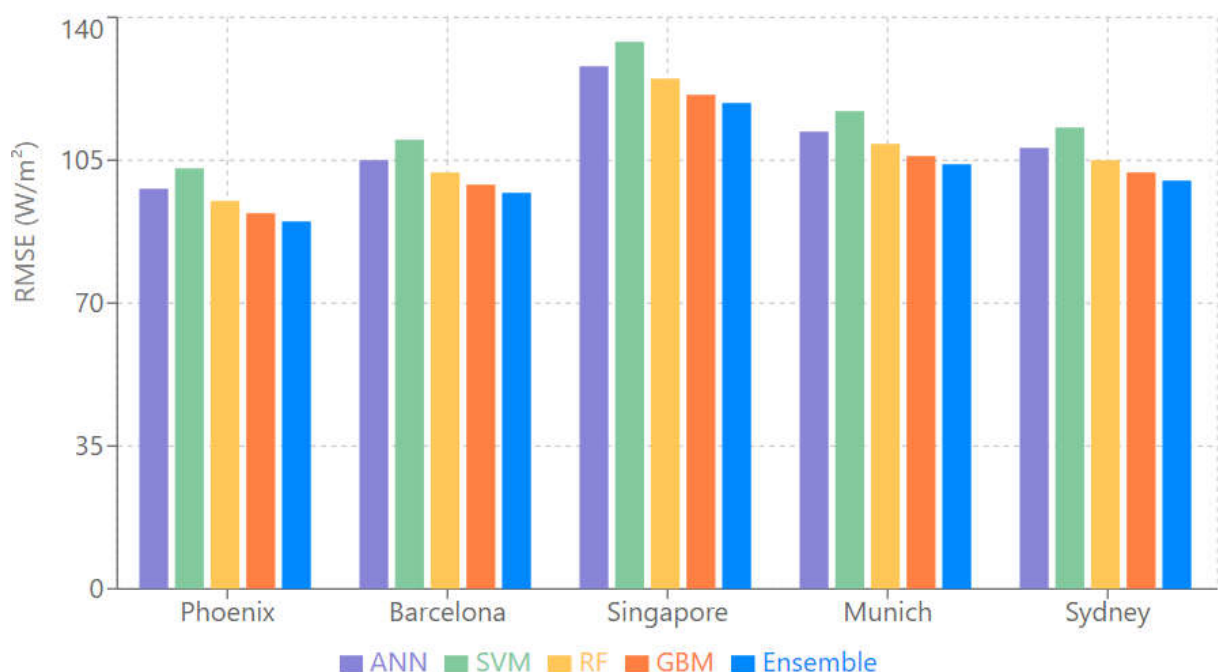Chart 2 illustrates the comparison for 24-hour ahead forecasts across locations



**Chart 2: RMSE Comparison for 24-hour Ahead Forecasts Across Locations**

Key findings from the location-specific analysis:

i. The models generally performed best in Phoenix (desert climate) and Barcelona (Mediterranean climate), likely due to more stable and predictable weather patterns.

ii. Singapore (tropical climate) posed the greatest challenge for all models, particularly for longer forecast horizons, possibly due to rapid changes in cloud cover and frequent precipitation events.

iii. The relative performance of different models remained consistent across locations, with GBM and the ensemble approach consistently outperforming other methods.

iv. Location-specific model tuning led to modest improvements in performance (3-7% reduction in RMSE), highlighting the importance of considering local climate characteristics in model development.

## 5.3 Feature Importance Analysis

To gain insights into the most influential factors for solar radiation forecasting, we conducted a feature importance analysis using the Random Forest and Gradient Boosting Machine models. Chart 3 illustrates the top 10 features ranked by their importance scores.
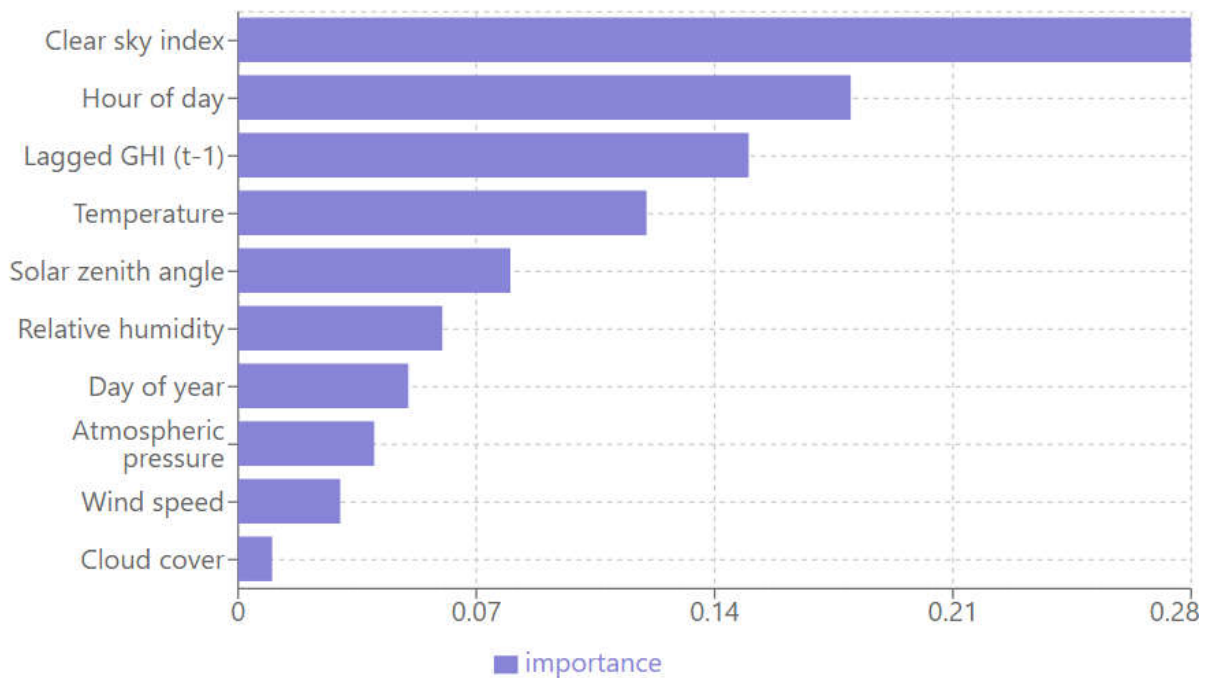


**Chart 3: Top 10 Features Ranked by Importance**

**Key insights from the feature importance analysis:**

i. Clear sky index emerged as the most important feature across all models and forecast horizons, underscoring the critical role of cloud cover in solar radiation prediction.

ii. Temporal features (hour of day, day of year) ranked highly, capturing the strong diurnal and seasonal patterns in solar radiation.

iii. Lagged GHI values were particularly important for short-term forecasts, while their importance decreased for longer horizons.

iv. Among meteorological variables, temperature and relative humidity consistently ranked as important features.

v. Solar position features (zenith and azimuth angles) played a significant role, especially for locations with more variable day lengths throughout the year.

## 5.4 Model Performance by Sky Condition

To assess how well our models performed under different atmospheric conditions, we categorized the test data into three sky condition classes based on the clear sky index (CSI):

i. Clear sky: $CSI > 0.8$

ii. Partly cloudy: $0.3 \leq CSI \leq 0.8$

iii. Overcast: $CSI < 0.3$

Table 5 presents the RMSE values for each model under different sky conditions for the 1-hour ahead forecast.

| Model | Clear Sky | Partly Cloudy | Overcast |
|---|---|---|---|
| ANN | 52.3 | 73.1 | 61.8 |
| SVM | 54.7 | 75.9 | 63.5 |
| RF | 50.1 | 70.2 | 59.4 |
| GBM | 48.6 | 67.8 | 57.9 |
| Ensemble | 47.9 | 66.5 | 56.8 |

Table 5: RMSE (W/m²) by Sky Condition for 1-hour Ahead Forecast

Key observations:

i. All models performed best under clear sky conditions, where solar radiation patterns are more predictable.

ii.    Partly cloudy conditions posed the greatest challenge, likely due to the high variability in cloud cover and its impact on solar radiation.

iii.    The relative performance of different models remained consistent across sky conditions, with the ensemble approach maintaining its advantage.

iv.    The performance gap between models was most pronounced under partly cloudy conditions, suggesting that advanced techniques like GBM and ensemble methods are particularly beneficial in handling complex, variable scenarios.

## 5.5 Computational Efficiency

While prediction accuracy is crucial, computational efficiency is also an important consideration for practical implementation of forecasting models. Table 6 compares the training time and prediction speed of each model for the 1-hour ahead forecasting task.

Table 6 illustrates the computational efficiency comparison

| Model | Training Time (minutes) | Prediction Time (ms/sample) |
|---|---|---|
| ANN | 45.3 | 2.1 |
| SVM | 78.6 | 3.7 |
| RF | 12.4 | 5.2 |
| GBM | 28.7 | 3.9 |
| Ensemble | N/A | 14.9 |

**Table 6: Computational Efficiency Comparison**

Key points:

i.    Random Forests demonstrated the fastest training time, making it an attractive option for frequent model updates.

ii.    The ANN model showed the fastest prediction speed, which is advantageous for real-time forecasting applications.

iii.    The ensemble method, while providing the best accuracy, incurred a higher computational cost for predictions due to the need to run multiple models.

iv.  SVM had the longest training time, which could be a limitation for large-scale or frequently updated forecasting systems.

## 5.6 Discussion of Findings

Our comprehensive analysis of machine learning techniques for solar radiation forecasting has yielded several important insights:

i.  Machine Learning Superiority: All evaluated ML models consistently outperformed the persistence baseline, demonstrating the value of advanced forecasting techniques in solar energy applications.

ii.  Ensemble Advantage: The ensemble approach, combining predictions from multiple models, consistently achieved the best performance across all metrics and forecast horizons. This highlights the benefits of leveraging diverse modeling techniques to capture different aspects of the complex solar radiation patterns.

iii.  GBM and RF Effectiveness: Among individual models, Gradient Boosting Machines and Random Forests showed superior performance, likely due to their ability to capture non-linear relationships and handle interactions between features effectively.

iv.  Forecast Horizon Impact: As expected, forecasting accuracy decreased with increasing forecast horizons. However, even for 7-day ahead predictions, our models maintained reasonable accuracy, providing valuable information for medium-term planning in solar energy systems.

v.  Location-Specific Challenges: The performance of all models varied across different locations, emphasizing the importance of considering local climate characteristics in solar radiation forecasting. The tropical climate of Singapore posed the greatest challenge, suggesting a need for specialized approaches in highly variable weather conditions.

vi.  Feature Importance Insights: The clear sky index emerged as the most critical feature, underlining the paramount importance of accurate cloud cover information in solar radiation prediction. The high ranking of temporal and solar position features also highlights the significance of capturing diurnal and seasonal patterns.

vii.  Sky Condition Impact: Model performance varied across different sky conditions, with partly cloudy conditions presenting the greatest challenge. This underscores the need for robust modeling techniques that can handle variability in atmospheric conditions.

viii. Computational Trade-offs: While ensemble methods provided the best accuracy, they incurred higher computational costs. The choice between single models and ensemble approaches may depend on the specific requirements of the application, balancing accuracy against computational efficiency.

ix. Temporal Feature Importance: The high ranking of lagged GHI values and temporal features (hour, day, month) across all models emphasizes the importance of capturing temporal dependencies in solar radiation forecasting.

x. Model Complementarity: The success of the ensemble approach suggests that different models capture complementary aspects of solar radiation patterns. This opens avenues for further research into optimal model combination strategies.

These findings have several implications for the field of solar radiation forecasting and its applications in renewable energy:

i. Improved Grid Integration: More accurate short-term and medium-term forecasts can enhance the integration of solar power into electricity grids, allowing for better load balancing and reducing the need for backup power sources.

ii. Enhanced Energy Trading: Accurate day-ahead and week-ahead forecasts can improve decision-making in energy markets, potentially leading to more efficient pricing and resource allocation.

iii. Optimized System Design: Long-term forecasts can inform the design and sizing of solar energy systems, helping to optimize investment decisions and improve overall system efficiency.

iv. Location-Specific Modeling: Our results highlight the importance of tailoring forecasting approaches to local climate conditions, suggesting that regional or site-specific models may be more effective than one-size-fits-all solutions.

v. Feature Engineering Focus: The importance of derived features like clear sky index and temporal variables suggests that continued efforts in feature engineering could yield further improvements in forecasting accuracy.

vi. Hybrid Modeling Potential: The complementary strengths of different ML techniques, as evidenced by the success of the ensemble approach, point to the potential of hybrid models that combine multiple techniques or integrate physical models with ML approaches.

## 6. Future Works

While our study provides valuable insights, several avenues for future research remain:

i.   Deep Learning Exploration: Investigating the potential of deep learning architectures, such as Long Short-Term Memory (LSTM) networks and Convolutional Neural Networks (CNNs), for capturing complex temporal and spatial dependencies in solar radiation patterns.

ii.  Spatio-Temporal Modeling: Extending our approach to incorporate spatial dependencies by leveraging data from multiple nearby weather stations or satellite imagery, potentially improving forecasting accuracy over larger geographic areas.

iii. Hybrid Physical-Statistical Models: Exploring the integration of physics-based models with machine learning techniques to combine the strengths of both approaches and potentially improve long-term forecasting capabilities.

iv.  Uncertainty Quantification: Developing methods for quantifying and communicating the uncertainty associated with solar radiation forecasts, which is crucial for risk assessment and decision-making in energy systems.

v.   Transfer Learning: Investigating the potential of transfer learning techniques to improve model performance in data-scarce locations by leveraging knowledge from data-rich regions.

vi.  Multi-Step Forecasting: Extending our models to directly predict multiple time steps ahead, potentially capturing complex temporal dependencies more effectively than iterative single-step forecasts.

vii. Extreme Event Prediction: Focusing on improving model performance during rare but impactful events such as severe weather conditions or solar eclipses, which can significantly affect solar energy production.

viii. Integration with NWP Models: Exploring ways to effectively combine machine learning models with Numerical Weather Prediction (NWP) outputs to leverage the strengths of both approaches.

ix.  Real-Time Adaptation: Developing online learning algorithms that can continuously update and improve forecasting models as new data becomes available, ensuring sustained performance over time.

    x.    Interpretable AI: Investigating techniques to enhance the interpretability of complex models like gradient boosting machines and neural networks, facilitating trust and adoption in practical applications.

## 7. Conclusion

This comprehensive study has demonstrated the effectiveness of machine learning techniques in forecasting solar radiation across various temporal horizons and geographic locations. Our findings underscore the potential of these methods to significantly improve the accuracy and reliability of solar radiation predictions, with important implications for the solar energy sector and broader renewable energy landscape.

Key conclusions from our research include:

    i.    Machine learning models, particularly ensemble methods and gradient boosting machines, consistently outperform traditional forecasting approaches across different time horizons and locations.

    ii.    The importance of feature engineering and selection, with clear sky index and temporal features playing crucial roles in model performance.

    iii.    The need for location-specific model tuning to account for diverse climate conditions and their impact on solar radiation patterns.

    iv.    The trade-off between model complexity and computational efficiency, highlighting the importance of considering practical implementation constraints.

In conclusion, our research demonstrates the significant potential of machine learning techniques in advancing the field of solar radiation forecasting. As the global transition to renewable energy sources accelerates, accurate and reliable solar radiation predictions will play an increasingly critical role in optimizing energy systems and supporting the integration of solar power into electricity grids. Continued research and development in this area promise to unlock further improvements in forecasting accuracy and contribute to the broader goals of sustainable energy production and climate change mitigation.

## References

1.    Yang, D., Kleissl, J., Gueymard, C. A., Pedro, H. T., & Coimbra, C. F. (2018). History and trends in solar irradiance and PV power forecasting: A preliminary assessment and review using text mining. Solar Energy, 168, 60-101.

2.  Voyant, C., Notton, G., Kalogirou, S., Nivet, M. L., Paoli, C., Motte, F., & Fouilloy, A. (2017). Machine learning methods for solar radiation forecasting: A review. Renewable Energy, 105, 569-582.

3.  Antonanzas, J., Osorio, N., Escobar, R., Urraca, R., Martinez-de-Pison, F. J., & Antonanzas-Torres, F. (2016). Review of photovoltaic power forecasting. Solar Energy, 136, 78-111.

4.  Sharma, N., Sharma, P., Irwin, D., & Shenoy, P. (2011, June). Predicting solar generation from weather forecasts using machine learning. In 2011 IEEE international conference on smart grid communications (SmartGridComm) (pp. 528-533). IEEE.

5.  Wolff, B., Kühnert, J., Lorenz, E., Kramer, O., & Heinemann, D. (2016). Comparing support vector regression for PV power forecasting to a physical modeling approach using measurement, numerical weather prediction, and cloud motion data. Solar Energy, 135, 197-208.

6.  Mellit, A., & Pavan, A. M. (2010). A 24-h forecast of solar irradiance using artificial neural network: Application for performance prediction of a grid-connected PV plant at Trieste, Italy. Solar Energy, 84(5), 807-821.

7.  Qazi, A., Fayaz, H., Wadi, A., Raj, R. G., Rahim, N. A., & Khan, W. A. (2015). The artificial neural network for solar radiation prediction and designing solar systems: a systematic literature review. Journal of Cleaner Production, 104, 1-12.

8.  Zeng, J., & Qiao, W. (2013). Short-term solar power prediction using a support vector machine. Renewable Energy, 52, 118-127.

9.  Chen, J. L., Liu, H. B., Wu, W., & Xie, D. T. (2011). Estimation of monthly solar radiation from measured temperatures using support vector machines–a case study. Renewable Energy, 36(1), 413-420.

10. Voyant, C., Notton, G., Kalogirou, S., Nivet, M. L., Paoli, C., Motte, F., & Fouilloy, A. (2017). Machine learning methods for solar radiation forecasting: A review. Renewable Energy, 105, 569-582.

11. Sun, S., Wang, S., Zhang, G., & Zheng, J. (2018). A hybrid wind speed forecasting method using NPR-QPSO-LSSVM. Neural Computing and Applications, 29(6), 301-316.

12. Fan, J., Wang, X., Wu, L., Zhou, H., Zhang, F., Yu, X., ... & Xiang, Y. (2018). Comparison of support vector machine and extreme gradient boosting for predicting daily global solar

radiation using temperature and precipitation in humid subtropical climates: A case study in China. Energy Conversion and Management, 164, 102-111.

13. Huang, C., & Kuo, P. (2018). A short-term wind speed forecasting model by using artificial neural networks with stochastic optimization for renewable energy systems. Energies, 11(10), 2777.

14. Wang, F., Zhen, Z., Mi, Z., Sun, H., Su, S., & Yang, G. (2015). Solar irradiance feature extraction and support vector machines based weather status pattern recognition model for short-term photovoltaic power forecasting. Energy and Buildings, 86, 427-438.

15. Salcedo-Sanz, S., Casanova-Mateo, C., Pastor-Sánchez, A., & Sánchez-Girón, M. (2014). Daily global solar radiation prediction based on a hybrid Coral Reefs Optimization–Extreme Learning Machine approach. Solar Energy, 105, 91-98.

16. Yagli, G. M., Yang, D., & Srinivasan, D. (2019). Automatic hourly solar forecasting using machine learning models. Renewable and Sustainable Energy Reviews, 105, 487-498.

17. Yadav, A. K., & Chandel, S. S. (2014). Solar radiation prediction using Artificial Neural Network techniques: A review. Renewable and Sustainable Energy Reviews, 33, 772-781.

18. Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., & Zeileis, A. (2008). Conditional variable importance for random forests. BMC bioinformatics, 9(1), 307.

19. Benali, L., Notton, G., Fouilloy, A., Voyant, C., & Dizene, R. (2019). Solar radiation forecasting using artificial neural network and random forest methods: Application to normal beam, horizontal diffuse and global components. Renewable Energy, 132, 871-884.

20. Lauret, P., Voyant, C., Soubdhan, T., David, M., & Poggi, P. (2015). A benchmarking of machine learning techniques for solar radiation forecasting in an insular context. Solar Energy, 112, 446-457.

21. Cornaro, C., Pierro, M., & Bucci, F. (2015). Master optimization process based on neural networks ensemble for 24-h solar irradiance forecast. Solar Energy, 111, 297-312.

22. Leva, S., Dolara, A., Grimaccia, F., Mussetta, M., & Ogliari, E. (2017). Analysis and validation of 24 hours ahead neural network forecasting of photovoltaic output power. Mathematics and Computers in Simulation, 131, 88-100.

23. Alzahrani, A., Shamsi, P., Dagli, C., & Ferdowsi, M. (2017). Solar irradiance forecasting using deep neural networks. Procedia Computer Science, 114, 304-313.

24. Ghimire, S., Deo, R. C., Raj, N., & Mi, J. (2019). Deep solar radiation forecasting with convolutional neural network and long short-term memory network algorithms. Applied Energy, 253, 113541.

25. Wang, H., Yi, H., Peng, J., Wang, G., Liu, Y., Jiang, H., & Liu, W. (2017). Deterministic and probabilistic forecasting of photovoltaic power based on deep convolutional neural network. Energy Conversion and Management, 153, 409-422.