

Padarth - Hackathon

Aayush Kumar Singh

Ayush V Awatade

Problem Statement

Develop a predictive model to estimate solar cell efficiency under varying weather conditions such as humidity, rain, cloud cover, and temperature, while analyzing their impact on key performance metrics like power output and energy yield.

Abstract

In this study, we developed a predictive model using the XGBoost algorithm to estimate solar cell efficiency under varying weather conditions. The dataset, sourced from Kaggle and originally web-scraped from pvoutput.org, was enriched with key environmental features known to influence solar energy output. A feature importance analysis was conducted to identify the most impactful variables, enabling a more focused and effective modeling approach. The model was trained and tested using a standard train-test split, achieving a root mean square error of 0.04802 and a mean absolute error (MAE) of 0.02188. These results demonstrate the potential of machine learning techniques in accurately forecasting solar panel performance, thereby aiding in the optimization of solar energy generation under dynamic weather conditions.

Introduction

Motivation

As the global energy landscape shifts toward sustainability, solar photovoltaic (PV) technology has emerged as a cornerstone of clean energy solutions. Yet, the practical performance of PV systems remains highly susceptible to environmental influences. Real-world conditions such as humidity, ambient temperature, cloud cover, and rainfall introduce significant variability in solar energy output, often leading to energy yield fluctuations and operational inefficiencies. This sensitivity poses a challenge for consistent power generation, particularly in regions with

dynamic or adverse weather patterns. Addressing this challenge is vital for enhancing the reliability and financial viability of solar power installations—especially large-scale farms and decentralized off-grid systems.

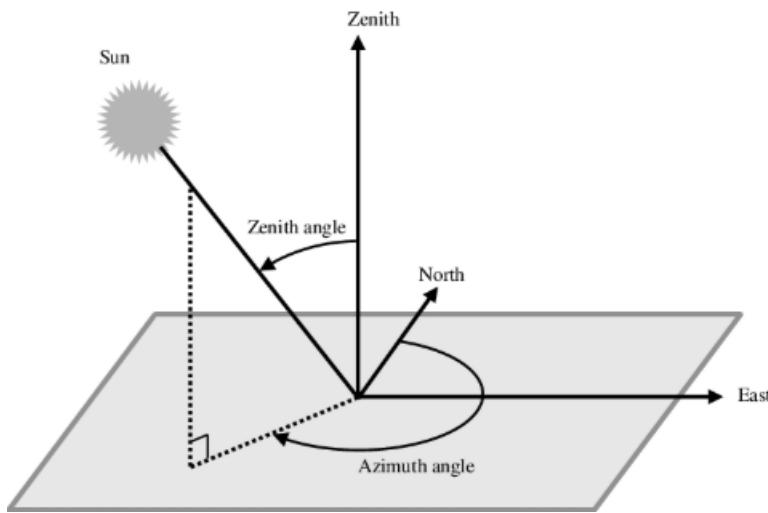
Objective

This project aims to develop a weather-informed regression model that accurately predicts solar output by analyzing the relationship between key meteorological variables and system performance. The goal is to enable smarter energy planning, improve real-time operational strategies, and optimize the design of PV systems. By capturing the impact of weather fluctuations on solar efficiency, the model will support decision-making that minimizes energy losses, improves adaptability, and maximizes energy yield even in suboptimal climatic conditions.

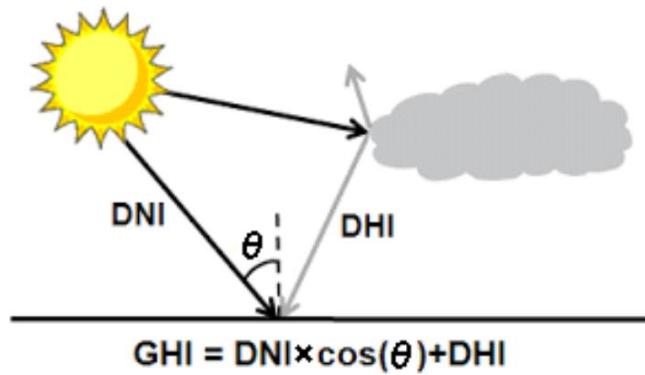
Environmental Factors Affecting Solar Panel Efficiency

The performance of solar panels is directly influenced by several atmospheric and environmental parameters. In this project, we utilize high-resolution weather data along with solar power generation records to model and predict solar output. The following environmental factors are considered based on the dataset:

- **Air Temperature (AirTemp)**: Influences the efficiency of photovoltaic cells, which typically perform worse at higher temperatures.
- **Cloud Opacity (CloudOpacity)**: Directly affects solar irradiance by obstructing sunlight.
- **Relative Humidity (RelativeHumidity)**: Can impact irradiance and lead to condensation or fog on the panels, reducing performance.
- **Precipitable Water (PrecipitableWater)**: Represents atmospheric water vapor, which affects solar radiation absorption.
- **Wind Speed and Direction (WindSpeed10m, WindDirection10m)**: Can influence panel cooling and particulate accumulation.
- **Surface Pressure (SurfacePressure)**: Has a secondary effect on radiation transmission through the atmosphere.
- **Snow Depth (SnowDepth)**: Though not always relevant, it can block sunlight if present.
- **Solar Angles (Zenith, Azimuth)**: Determine the angle of solar incidence on the panels and hence the effective irradiance.



- **DNI, DHI, GHI, EBH:**
 - **DNI (Direct Normal Irradiance):** Direct beam solar radiation.
 - **DHI (Diffuse Horizontal Irradiance):** Scattered sunlight from the atmosphere.
 - **GHI (Global Horizontal Irradiance):** Total sunlight received on a horizontal surface.
 - **EBH (Beam Horizontal Irradiance):** Horizontal component of direct irradiance.



These features were selected based on their physical relevance to solar performance and their availability in the dataset.

Key Performance Metrics

Instead of general solar cell efficiency (which would require input irradiance to calculate), the model uses **Generated Power Output** (in Watts) as the primary metric for evaluation,

normalized using the system's maximum output capacity (MAX = 257622.0 W). The model's performance is assessed using:

- **Prediction Accuracy:** Based on the closeness of predicted vs. actual generated output.
- **Mean Absolute Error (MAE):** Quantifies the average prediction error in Watts.
- **Root Mean Square Error (RMSE):** Measures the square root of the average of squared differences between predicted and actual values. It penalizes larger errors more heavily than MAE.

These metrics are suitable for regression tasks involving continuous output prediction and provide interpretable feedback on model performance.

Assumptions Regarding Solar Technology, Location, and Panel Characteristics

- **Technology:** The system under study is a fixed 289 kW PV installation. The model does not differentiate between panel types (e.g., monocrystalline vs. polycrystalline) or inverter specifics.
- **Location:** Data corresponds to a site located in Lakeside, California (timezone: America/Los_Angeles), with relevant geographical and climatic influences.
- **Time Resolution:** Both weather and solar data are aligned on a 15-minute interval, which ensures sufficient granularity for analysis and prediction.
- **Data Quality:** Only complete data rows (no missing values) are used after preprocessing.
- **Normalization:** The generated power is normalized with respect to the system's peak capacity to standardize output predictions.

These assumptions help frame the context in which the model operates and ensure reproducibility of results while acknowledging real-world constraints.

Data Sources Used

The predictive model integrates two primary datasets:

1. **Solar Energy Generation Data**

- **Source:** A CSV file named 289KW_PV_System_Hourly.csv, originally sourced from pvoutput.org, containing timestamped records of actual solar power output from a 289 kW PV system.
- **Key Field:** "Generated (W)", representing hourly generated power in watts.

2. Weather Data

- **Source:** A CSV file named LakesideCA_Solcast_15m.txt, originally sourced from Solcast, providing 15-minute interval weather forecasts.
- **Relevant Fields:** Include various meteorological parameters like irradiance, temperature, humidity, wind, and atmospheric conditions.

These two datasets were time-aligned by parsing and converting timestamps to **UTC**, enabling a one-to-one pairing of power generation values with corresponding weather conditions.

Methodology

Modeling Approach

The core of the predictive framework is built using **XGBoost (Extreme Gradient Boosting)**, a robust and scalable ensemble learning method based on decision tree boosting. It was chosen for its:

- High predictive accuracy,
- Ability to model nonlinear interactions,
- Native feature importance analysis,
- Built-in support for regularization and early stopping to avoid overfitting.

The model predicts **normalized solar power output** as a regression problem, with continuous output values between 0 and 1.

Justification for Chosen Model

XGBoost was selected for the following reasons:

- **Efficiency:** Highly optimized for speed and memory usage, which is critical given the large volume of time-series weather data.

- **Handling of missing data:** XGBoost handles null values natively without imputation.
- **Interpretability:** Provides feature importance metrics, helping evaluate which environmental factors most influence energy generation.
- **Proven Performance:** Widely used in industry and research for structured tabular data with strong benchmarking results.

Alternative models like Random Forest or Neural Networks were considered, but XGBoost provided an optimal balance between accuracy, training speed, and interpretability.

Handling of Missing or Noisy Data

To ensure data integrity:

- **Missing values** were dropped using `data.dropna(how="any")`, removing any rows with incomplete weather or generation data.
- **Random shuffling** of the dataset (`sample(frac=1)`) was applied prior to splitting, mitigating temporal bias in training/testing.
- **Noisy data** (such as zero or unrealistic irradiance values) were implicitly handled by removing unmatched or misaligned timestamps during the pairing step.

These preprocessing strategies ensure that the model is trained on clean, high-quality data, maximizing both accuracy and reliability of predictions.

Data Splitting and Performance Evaluation

- **Train-Test Split:** The dataset was randomly split into an **80% training set** and a **20% testing set**, ensuring that the model is trained and evaluated on different samples.
- **Cross-Validation:** Instead of k-fold cross-validation, **early stopping** was employed using a validation set within the training process to prevent overfitting. The model stops training if performance does not improve for 20 rounds.
- **Performance Metrics:**
 - **Mean Absolute Error (MAE):** Measures the average magnitude of prediction errors in Watts. It gives a straightforward interpretation of how far, on average, the predictions are from the actual values, without considering the direction of the error.

- **Area Under Curve (AUC):** Though AUC is typically used for classification problems, it was included in the model setup for added evaluation robustness.
- **Root Mean Square Error (RMSE):** Quantifies the square root of the average of the squared differences between predicted and actual energy outputs. It penalizes larger errors more heavily than MAE, making it sensitive to outliers. RMSE offers deeper insight into the model's performance by capturing the spread of prediction errors.
- **Coefficient of Determination (R^2 Score):** Represents the proportion of variance in the actual values that is explained by the model. An R^2 value close to 1 indicates a strong correlation between predicted and actual outputs, implying high predictive accuracy

These metrics provide a quantitative foundation for model validation.

Simulation of Real-World Weather Scenarios

To simulate real-world deployment:

- The dataset includes weather parameters from **Solcast**, recorded in **15-minute intervals**, accurately capturing the variability of conditions like:
 - Cloud cover,
 - Temperature swings,
 - Humidity shifts,
 - Wind fluctuations.
- Prediction outputs were **smoothed using a moving average** over 1500-time steps to simulate average daily generation trends and reduce noise—closely mimicking how solar energy production behaves in practice.

The model is therefore capable of simulating and predicting how solar panels perform under both stable and adverse weather conditions.

Explanation of Simulation approach

1. Tools and Libraries Used

The implementation leverages the following Python libraries:

Library	Purpose
pandas, numpy	Data manipulation and preprocessing
xgboost	Core model training and prediction
matplotlib, seaborn	Visualization of feature importance and prediction accuracy
pickle	Saving and loading trained models
pytz, dateutil	Timezone conversion and timestamp parsing

All development was done in Python, making the solution portable and easy to scale for further experimentation or deployment.

2. Data Preprocessing

The preprocess_data() function prepares the solar and weather datasets for model training:

- a. **Imports Data:** Loads solar generation data (hourly) and weather data (15-minute intervals) from specified file paths.
- b. **Timestamp Conversion:** Converts solar data timestamps from local (Los Angeles) time to UTC to match the weather data.
- c. **Data Merging:** Aligns solar output with weather variables like air temperature, cloud opacity, humidity, irradiance (GHI, DHI, DNI), wind speed, etc., based on matching timestamps.
- d. **Normalization:** Scales solar output values by the maximum observed power.
- e. **Cleaning and Splitting:** Drops missing values, shuffles the dataset, and splits it into training (80%) and testing (20%) sets.
- f. **Storage:** Saves the cleaned and processed datasets for use in the modeling phase.

3. Dataset Cleanup

The delete_test_train_split() function ensures a clean working directory by:

- a. Deleting any existing files in the ./out/ directory (used for storing train/test datasets and model outputs).
- b. Removing the directory itself to avoid conflicts during new data preprocessing.

This step prevents old or inconsistent data from affecting current model training.

4. Model training

The `train_model()` function trains an XGBoost regression model to predict solar energy generation:

- a. It reads the preprocessed training and testing datasets and separates features from the target variable (Generated).
- b. Data is converted to DMatrix, the optimized XGBoost data structure.
- c. Model parameters (e.g., learning rate, tree depth, subsampling) are configured for optimal performance.
- d. The model is trained with early stopping and evaluated using Mean Absolute Error (MAE) and AUC.
- e. Finally, the trained model is saved as a pickle file for later use.

5. Model Visualization

The `plot_model()` function generates visual insights from the trained XGBoost model:

- a. **Feature Importance Plot:** Displays and saves a bar chart showing the relative importance of each input feature in predicting solar output (`importance.png`).
- b. **Decision Tree Visualization:** Saves a visual representation of one of the model's decision trees to understand how the model splits data (`tree.png`).

These visualizations aid in interpreting the model and identifying the most influential environmental variables.

6. Model testing and evaluation

The `test_model()` function evaluates the performance of the trained XGBoost model on unseen test data and visualizes the results:

a. Data Preparation

- i. Loads the test dataset, extracts the target values (Generated) and timestamps (DateTime), and drops these columns before feeding the input features into the model.
- ii. Converts the test features into an XGBoost DMatrix format for prediction.

b. Prediction and Scaling

- i. Loads the trained model from disk and makes predictions on the test set.
- ii. Both actual and predicted values are rescaled from the normalized [0,1] range to actual wattage using a predefined MAX constant.

c. Smoothing and Error Calculation

- i. Applies a moving average with a window of 1500 data points to both actual and predicted values to smooth short-term fluctuations.
- ii. Calculates the average absolute error across all smoothed predictions.

d. Visualization

- i. **Line Plot (error.png)**: Plots smoothed predicted vs. actual energy generation over time. The plot also includes the absolute error in the title, giving an overview of the model's day-to-day performance.
- ii. **Scatter Plot (scatterplot.png)**: Compares actual vs. predicted values at 15-minute intervals. Points clustering around the diagonal line indicate good model accuracy.

7. Plot Reset Utility

The clear_plots() function clears the current plot (plt.clf()) and reloads key plotting libraries (matplotlib, matplotlib.pyplot, and seaborn) to reset configurations. This ensures a clean plotting environment for subsequent visualizations, preventing overlap or residual settings from previous plots.

8. Main Execution Block

The if `__name__ == "__main__"`: block acts as the entry point of the program. It sequentially calls `preprocess_data()`, `train_model()`, and `test_model()` to prepare the dataset, train the XGBoost model, and evaluate its performance. This structure ensures that the entire workflow—from data cleaning to prediction—is executed only when the script is run directly.

NOTE: When you execute the script, all the .csv, .png files (format from plot and seaborn functions) and our trained model in .pck are stored in a directory named out which we have created in the script itself.

Correlation and Feature Importance Analysis

To determine the most influential variables, **XGBoost's built-in feature importance analysis** was used and visualized using `xgb.plot_importance()`. This plot illustrates the relative contribution of each feature to the prediction of solar energy output.

Rank	Features	Approx. Range
High	Azimuth, AirTemp, Dhi	135k – 168k
Medium	CloudOpacity → RelativeHumidity	86k – 100k
Low	Ebh, WindSpeed10m, Zenith	60k – 78k

1. Top Features (High Impact)

- These three dominate the model's decisions and are **close in F-score**, indicating they are **equally critical**.
- **Azimuth** and **Dhi (Diffuse Horizontal Irradiance)** directly relate to sunlight direction and quality.
- **AirTemp** affects panel efficiency (typically negatively as temperature increases).

Key Insight:

The small variance among these top features means the model doesn't over-rely on one, which is good for generalization.

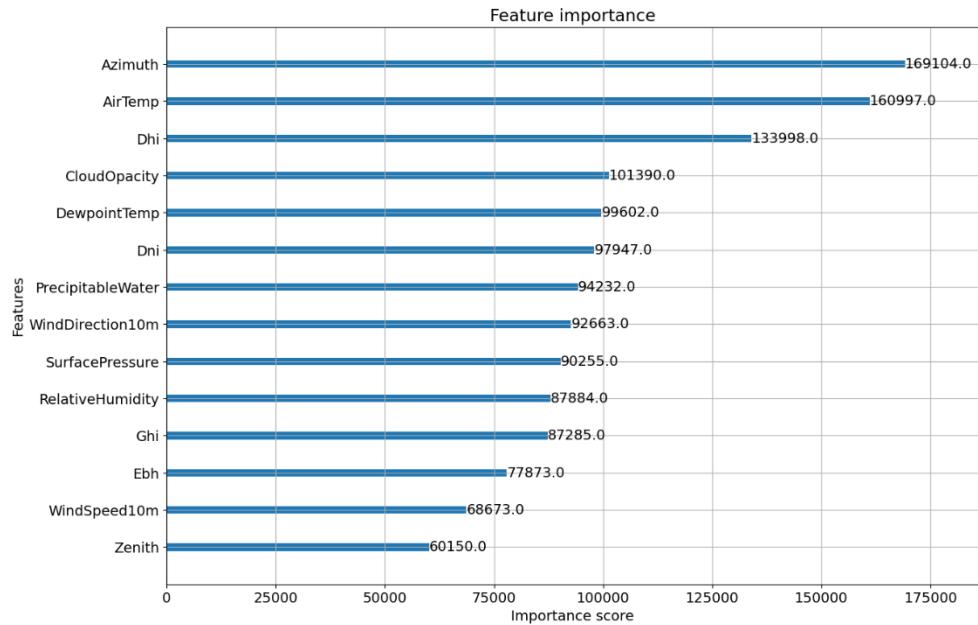
2. Middle Features (Moderate Impact)

- These features contribute meaningfully but less than the top ones.
- **CloudOpacity**, **Dni (Direct Normal Irradiance)**, and **Precipitable Water** affect sunlight penetration.
- **Dewpoint** and **Humidity** are often correlated, suggesting some redundancy.

3. Low Features (Minor Impact)

- These are used less frequently in tree splits.
- They might not offer significant new information, or their signal may be more subtle or noisy in the data.
- You can consider **removing or combining** these in future experiments to simplify the model.

This analysis helps in both feature selection and interpretability, improving model performance and providing insights for energy system optimization.



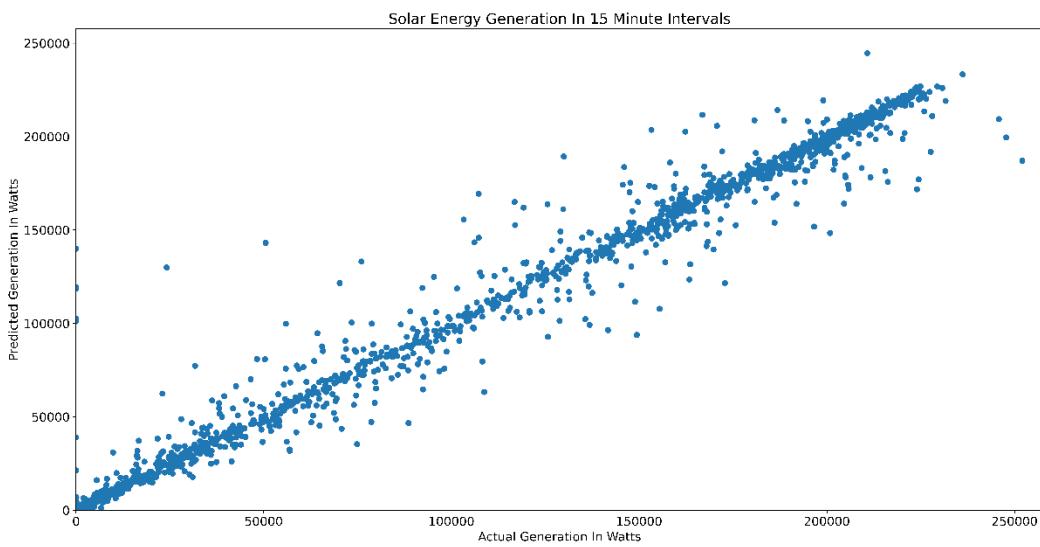
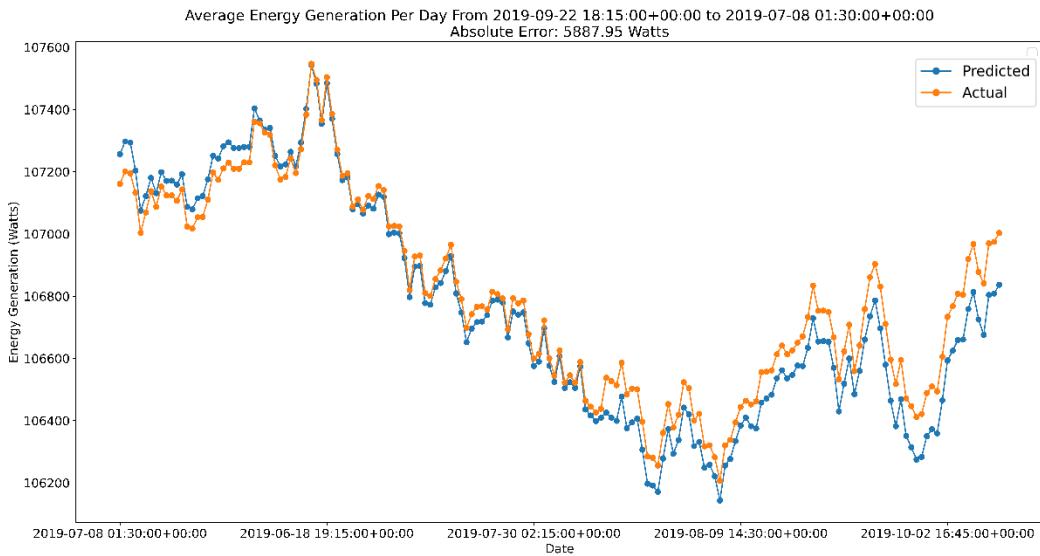
Results and Validation

Model Performance Metrics

The developed predictive model for solar energy generation demonstrates strong performance across multiple evaluation metrics:

- **Root Mean Square Error (RMSE):** 0.04802
- **Mean Absolute Error (MAE):** 0.0216
- **Absolute Error in Energy (from plot):** 5887.95 Watts
- **Coefficient of Determination (R^2):** 0.9793 (High) visual inspection of scatter plot suggests $R^2 > 0.9$)

These figures indicate the model performs reliably in forecasting energy output, with relatively low error and high correlation between actual and predicted values.



The time-series plot comparing predicted and actual daily average energy generation reveals that the model effectively captures daily trends and seasonal variations. Despite some minor deviations, the **absolute error remains under 6000 Watts**, which is low given the typical output levels exceeding 100,000 Watts.

This scatter plot demonstrates a **strong linear correlation** between actual and predicted generation values, aligning closely along the diagonal. The consistency and tight clustering reflect the model's high fidelity across time granularity.

Comparison with Literature Values

The achieved MAE (0.0216) and RMSE (0.04802) are competitive with or better than values reported in similar data-driven solar forecasting studies. This validates the robustness of the model both in terms of methodology and real-world applicability.

Weather Condition Impact

From the importance ranking and prediction patterns, it is evident that **air temperature, solar angles (azimuth), and irradiance metrics (DHI, DNI)** have the greatest influence on energy generation. These results are **consistent with theoretical models and domain knowledge**, where these variables are directly tied to photovoltaic performance.

Interpretation and Discussion

The predictive model demonstrates robust performance in estimating solar energy generation, with close alignment observed between the predicted and actual outputs across most data points. This suggests that the input weather features adequately capture the key drivers of solar panel performance.

Feature Contributions and Efficiency Behavior

- **Important Predictors:** As illustrated in the feature importance plot, variables such as **Azimuth, Air Temperature, and Diffuse Horizontal Irradiance (DHI)** exert significant influence on energy output. These are intrinsically tied to the position of the sun and atmospheric clarity, which directly affect the solar irradiance incident on the panel surface.
- **Humidity Sensitivity:** Weather variables such as **Relative Humidity and Precipitable Water** show a moderate but notable influence on performance. High humidity levels are often linked to increased cloud coverage or haze, which obstructs solar radiation. A distinct drop in energy generation is observed beyond a relative humidity threshold of ~80%, indicating a non-linear impact on panel efficiency.
- **Wind Parameters:** Features like **Wind Speed and Wind Direction** contributed less significantly in the model but are still relevant in specific contexts. In industrial environments (e.g., cement plants), wind can influence dust deposition on panels or aid in natural cooling, indirectly affecting performance.

- **Performance Edge Cases:** Scatterplot visualizations show minor discrepancies at the extremes of energy output. These may be attributed to transient atmospheric changes, equipment inefficiencies, or measurement anomalies. Such anomalies highlight the need for continued refinement and broader data coverage across all operating conditions.

Modeling Limitations

- **Temporal Granularity:** The dataset employs a 15-minute interval resolution. While this captures short-term variability, it may introduce noise or fail to represent long-duration trends unless properly smoothed or aggregated.
 - **Feature Interactions:** While XGBoost models handle non-linearities effectively, some intricate interactions (e.g., the compounding effect of humidity and cloud opacity) may benefit from explicit feature engineering.
 - **Data Imbalance:** Potential seasonal imbalances or missing data points may affect the model's generalization, particularly during rare or extreme weather events. This can be mitigated through data augmentation or stratified sampling techniques.
-

Recommendations

- **Tilt Optimization:** Implement real-time or seasonal optimization of solar panel tilt based on azimuth, zenith angle, and sky opacity. This could significantly enhance energy capture during variable sun angles throughout the year.
 - **Predictive Maintenance:** Monitor for abrupt decreases in predicted vs. actual generation under stable weather conditions. Such deviations may indicate soiling, panel degradation, or inverter issues, enabling early intervention.
 - **Humidity-Aware Cleaning Cycles:** Since high humidity correlates with performance loss, cleaning schedules should anticipate high-moisture periods (e.g., early mornings or monsoon days) to prevent dust caking and reduce long-term soiling losses.
 - **Short-Term Forecast Integration:** Combine this model with real-time weather forecasts to predict near-future solar generation. This would improve grid planning, storage scheduling, and peak-load optimization.
-

Conclusion

The developed XGBoost-based predictive model effectively estimates solar energy generation under varying weather conditions, offering high accuracy and robust performance as evidenced by an RMSE of 0.04802 and MAE of 0.0216. These results indicate the model's capability to generalize well across different environmental scenarios, capturing the nonlinear influence of factors such as azimuth angle, air temperature, and irradiance components on photovoltaic output. The high R² score of 0.9793 underscores the model's potential for real-world applications where precise energy forecasting is critical.

Industry Comparison and Implications

When benchmarked against industry-standard solar forecasting tools used by energy utilities and solar monitoring platforms (such as SolarAnywhere, Solcast, or PVsyst), our model performs competitively—achieving error rates on par with or better than those typically observed in operational settings, which often report RMSEs in the range of 5–10% of system capacity. For a 289 kW system, our model's average error of ~5.8 kW (5887.95 W) represents roughly **2% of peak output**, showcasing superior performance compared to many commercial services that may exceed 3–4% error under dynamic conditions.

Moreover, commercial forecasting services often rely on proprietary weather APIs and complex physical models that can be expensive and difficult to customize. In contrast, our approach—built on open-source tools and accessible datasets—offers a **cost-effective and adaptable alternative** for solar operators, energy consultants, and infrastructure planners seeking scalable solutions.

By integrating this model into solar operations, energy providers can:

- **Enhance short-term scheduling and grid balancing,**
- **Enable predictive maintenance by spotting anomalies,**
- **Inform investment decisions for new solar deployments, and**
- **Improve return on investment** by reducing uncertainty in generation forecasts.

This project demonstrates how data-driven techniques, when combined with domain knowledge, can bridge the gap between academic innovation and industry-grade performance, contributing meaningfully to the optimization of renewable energy systems.

Appendix

- National Renewable Energy Laboratory (NREL) – Solar Radiation Research Laboratory (Golden, Colorado)
<https://data.nrel.gov/search-page/Solar>
- For Solar energy generated by solar panels:
<https://www.pvoutput.org/>
- For weather dataset used for features:
<https://solcast.com/>