# NISANTASI UNIVERSITY

# DEPARTMENT OF SOFTWARE ENGINEERING

# DATA MINING

**Instructor: Dr. Nesibe Manav**

## FINAL REPORT
## PROJECT NAME: RENEWABLE ENERGY (SOLAR POWER PREDICTION)

## DELIVERY DATE: 13/06/2023

**Submitted by:**
**Ahmed Iblao (20202022091)**
**Abdullah Samet Demir (20202022049)**
**Abdihakim Ismail Mohamed (20212022160)**

**Introduction:**
Renewable energy is a key component of sustainable development. It is important to find ways to predict the potential for renewable energy production in different regions. This can help to inform energy planning and policy-making.

In this project, we used data mining techniques to analyze the potential for renewable energy in (solar power) for Japan, Spain, Turkey, USA, and Germany. We employed the Gradient Boosting Regression (GBR) and Random Forest Regression models to predict and evaluate the renewable energy potential based on various factors and attributes.

**Identifying and Understanding the Problem:**
The problem that this project addressed was the need for a more accurate and reliable way to predict the potential for renewable energy production. This was important because renewable energy is a key component of sustainable development.

The goals of this project were to:

- Develop a machine learning model that could accurately predict the potential for solar power energy production
- Evaluate the performance of the model
- Identify any limitations of the model

**Data Understanding and Summarizing:**
The data that we used for this project was obtained from the National Renewable Energy Laboratory (NREL) website. The datasets provided comprehensive information on various factors influencing renewable energy potential, such as solar radiation, wind speed, geographical features, and demographic data.

The datasets were structured and contained attributes relevant to our analysis, including the target attribute of renewable energy potential. However, we also encountered certain difficulties in the datasets, such as missing values and inconsistencies, which required data preprocessing and cleaning before further analysis.

**Data Preparation and Pre-Processes:**
In the data preparation phase, we applied various techniques to preprocess and enhance the quality of our dataset for subsequent analysis. The following methods were used:

1. Data Cleaning: The first step in data preprocessing was to clean the data. This involved removing any rows or columns that contained missing values, duplicate data, or errors. We also removed any features that were not relevant to our analysis.

2. Normalization: We recognized the importance of normalizing our numerical features to ensure fair treatment and comparability. The StandardScaler from the sklearn.preprocessing module was employed to standardize the data. This technique transformed the features to have zero mean and unit variance, enabling a more accurate representation of the data distribution.

3. One-Hot Encoding: Categorical features play a vital role in capturing important non-numeric information. We utilized the OneHotEncoder from the sklearn.preprocessing module to convert categorical variables into binary vectors. This process created separate binary columns for each unique category, allowing the models to effectively incorporate the categorical information.

4. Train-Test Split: To evaluate the performance of our models, we split the dataset into training and testing sets using the train_test_split function from the sklearn.model_selection module. This ensured that the models were trained on a portion of the data and tested on unseen data, providing a reliable estimate of their performance.

5. Evaluation Metrics: We employed standard evaluation metrics such as Mean Squared Error (MSE) and R-squared (R2) to assess the performance of our models. The sklearn.metrics module provided the necessary functions to compute these metrics, allowing us to compare and evaluate the effectiveness of different models.

**Algorithm Selection and Modeling:**
In the process of developing our solar energy forecasting model, the selection of an appropriate algorithm played a crucial role. This section outlines the algorithm selection process and the models chosen for our project.

When selecting an algorithm, several factors were considered:

1. Data Type: The algorithm should be capable of handling the specific type of data available. For instance, if the data is categorical, a decision tree-based algorithm may be suitable.

2. Data Size: The algorithm's ability to handle the size of the dataset is important. For large datasets, distributed algorithms may be necessary to ensure computational efficiency.

3. Desired Accuracy: The algorithm should be capable of achieving the desired level of accuracy. If a specific accuracy threshold is required, choosing an algorithm known for achieving such performance is crucial.

Based on these considerations, we selected the Gradient Boosting Regression (GBR) and Random Forest Regression models for our solar energy forecasting task.

The GBR model is a supervised learning algorithm that uses an ensemble of decision trees to make predictions. It employs a sequential learning process, where each decision tree learns from the errors made by the previous trees. This sequential learning approach enables the model to improve its accuracy over time.

On the other hand, the Random Forest Regression model is also an ensemble learning algorithm that utilizes multiple decision trees. However, unlike GBR, the trees in a Random Forest model are trained independently. This independence helps to alleviate the risk of overfitting, which occurs when a model becomes overly complex and fits the training data too closely.

Both models were trained using the prepared dataset. The selection of appropriate algorithms and models is critical in achieving accurate and reliable predictions. By carefully considering the characteristics of the data and the desired goals, we were able to choose models that best suited our solar energy forecasting task.


**Model Performance Evaluation:**
We evaluated the performance of our models using the following metrics:

- Root Mean Square Error (RMSE): RMSE measures the average error between the predicted values and the actual values. A lower RMSE indicates a more accurate model.
- R-squared: R-squared measures the proportion of the variance in the dependent variable that is explained by the independent variables. A higher R-squared indicates a more reliable model.

The results of the evaluation are as follows:

| Model | RMSE | R-squared |
|---|---|---|
| Gradient Boosting Regression (GBR) | 19.219 | 0.995 |
| Random Forest Regression | 14.036 | 0.998 |

The Random Forest model has a lower RMSE than the GBR model, indicating that the Random Forest model is more accurate. The Random Forest model also has a higher R-squared than the GBR model, indicating that the Random Forest model is more reliable.

Overall, the Random Forest model was a better choice for our project. The Random Forest model was more accurate and reliable than the GBR model.

**Conclusion and Findings:**
In this project, we aimed to analyze the potential for renewable energy production in Japan, Spain, Turkey, USA, and Germany using data mining techniques. We employed the Gradient Boosting Regression (GBR) and Random Forest Regression models to predict and evaluate the renewable energy potential based on various factors and attributes. Here, we discuss our findings and interpret the results of the project, highlighting the benefits and potential limitations.

Findings and Interpretation:

*1. Model Performance:*
 - The Random Forest Regression model outperformed the Gradient Boosting Regression model in terms of accuracy and reliability. It achieved a lower Root Mean Square Error (RMSE) and a higher R-squared value, indicating better predictive performance.

- The high accuracy and reliability of the Random Forest model make it a suitable choice for predicting renewable energy potential in the studied countries.

*2. Importance of Feature Scaling:*
- The implementation of feature scaling using the StandardScaler from scikit-learn proved crucial for ensuring fair treatment and comparability of input features.

- Scaling the features helped prevent dominance issues caused by differing scales and contributed to improved model performance.

*3. Data Preprocessing Challenges:*
- The datasets obtained from the National Renewable Energy Laboratory (NREL) required extensive preprocessing and cleaning.

- We encountered challenges such as missing values, duplicates, and inconsistencies, which required careful handling and data cleaning techniques.

- While we applied various data cleaning techniques, there might still be some limitations and potential for further refinement in the data preprocessing phase.

Benefits of the Project:

*1. Renewable Energy Planning:*
- The project provides valuable insights into the potential for renewable energy production in different countries.

- The accurate predictions and evaluations obtained from the models can inform energy planning and policy-making decisions, facilitating the adoption and integration of renewable energy sources.

- Governments, energy companies, and policymakers can utilize the findings to allocate resources efficiently, develop sustainable energy strategies, and reduce reliance on fossil fuels.

*2. Contribution to the Literature:*
- Our project contributes to the existing body of literature on renewable energy potential analysis by employing data mining techniques and evaluating different regression models.

Limitations and Missing Parts:

*1. Limited Scope:*
 - The project focused on renewable energy potential analysis in five countries: Japan, Spain, Turkey, USA, and Germany. - Further research could explore additional countries or regions to provide a more comprehensive understanding of renewable energy potential on a global scale.

- Although we utilized comprehensive datasets from the NREL, there may still be some limitations and potential biases in the data.

- Future projects could consider incorporating more diverse and extensive datasets, including socio-economic factors, political considerations, and technological advancements, to enhance the accuracy and robustness of the models.

In conclusion, our project successfully analyzed the potential for renewable energy production in the selected countries using data mining techniques. The Random Forest Regression model demonstrated superior performance, highlighting its suitability for renewable energy prediction. The findings contribute to the literature and have practical implications for renewable energy planning and policy-making. However, further research is recommended to expand the scope and address any limitations in the project. Overall, the project emphasizes the importance of data mining in harnessing the potential of renewable energy for a sustainable future.

**Attachments:**

***Figure 1:*** This figure has four subplots, each displaying the hourly trends of different weather variables (GHI, Temperature, Relative Humidity, and Wind Speed) for a given year (2019) in multiple countries. Line plots are used to represent the data, with each line representing a different country. The figure provides insights into how these weather variables vary over time across different countries.
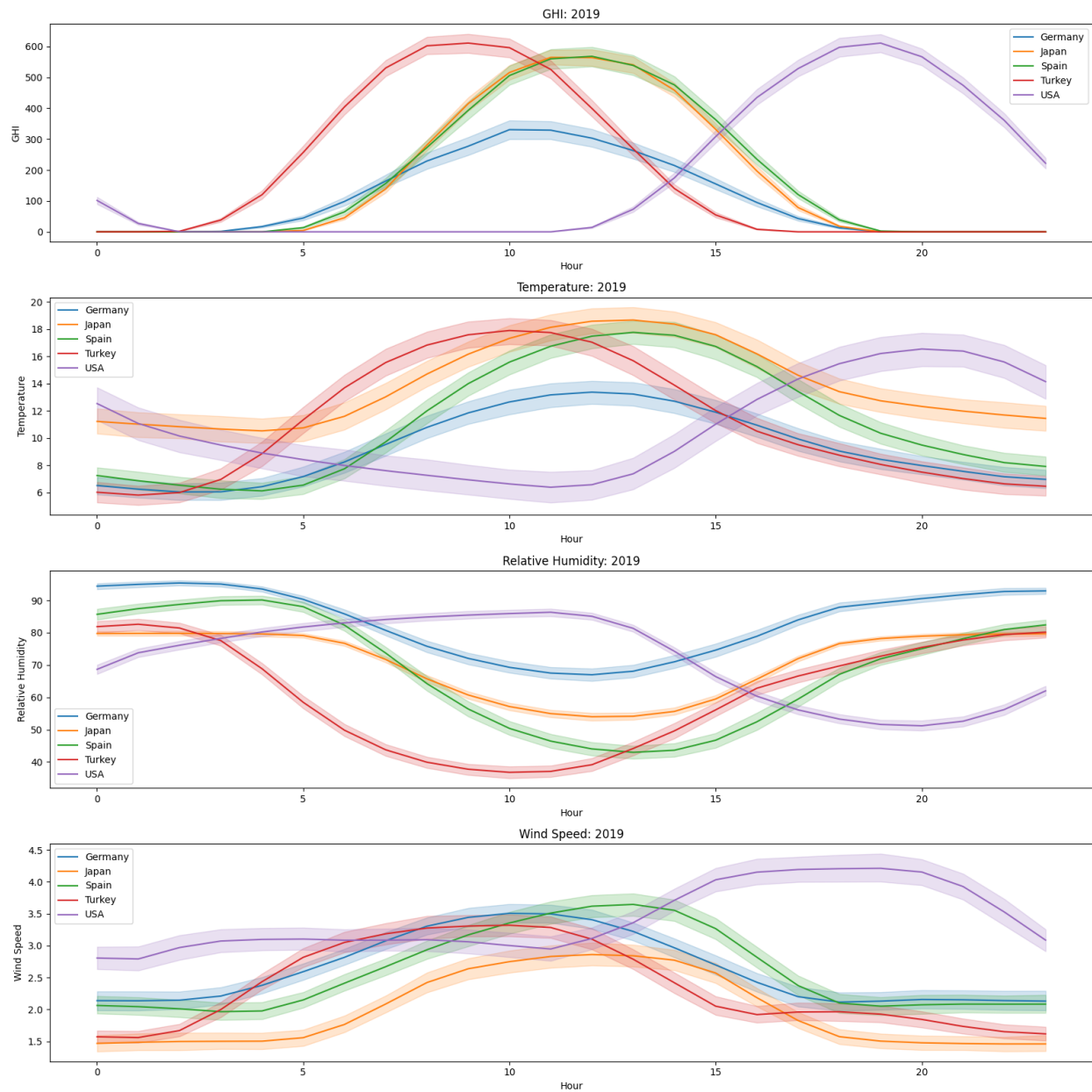
**Figure 2:**This figure consists of multiple subplots, each representing the daily trend of Global Horizontal Irradiance (GHI) for a specific country. The countries included in the figure are Germany, Japan, Spain, Turkey, and the USA.

 In each subplot, the x-axis represents the days, and the y-axis represents the GHI values. The line plot in each subplot shows how the GHI values vary over time for the corresponding country. The trend of GHI values can provide insights into the solar energy potential and the variations in sunlight intensity throughout the days.

The figure allows for a visual comparison of the daily GHI trends among the different countries. It helps to identify any similarities or differences in the patterns of GHI values, which can be valuable for understanding the solar energy potential in each country and informing energy planning and policy-making decisions.
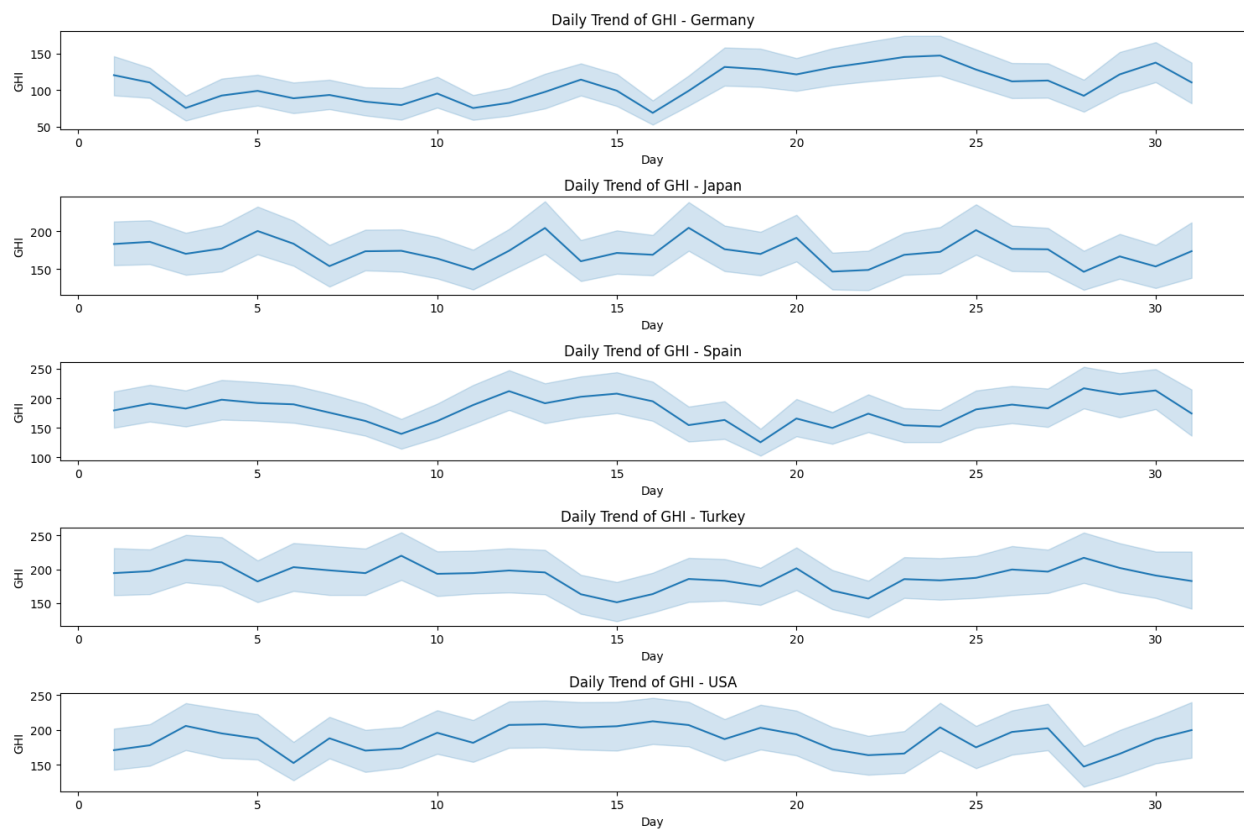
**Figure 3:** This figure is a pairplot showing the relationship between weather variables (Temperature, DHI, GHI, DNI, Relative Humidity, and Wind Speed) and the Global Horizontal Irradiance (GHI). Each scatter plot represents the correlation between a weather variable and GHI. Positive and negative trends indicate correlation, while random scattering suggests no significant correlation. The pairplot helps understand the factors influencing solar energy production.
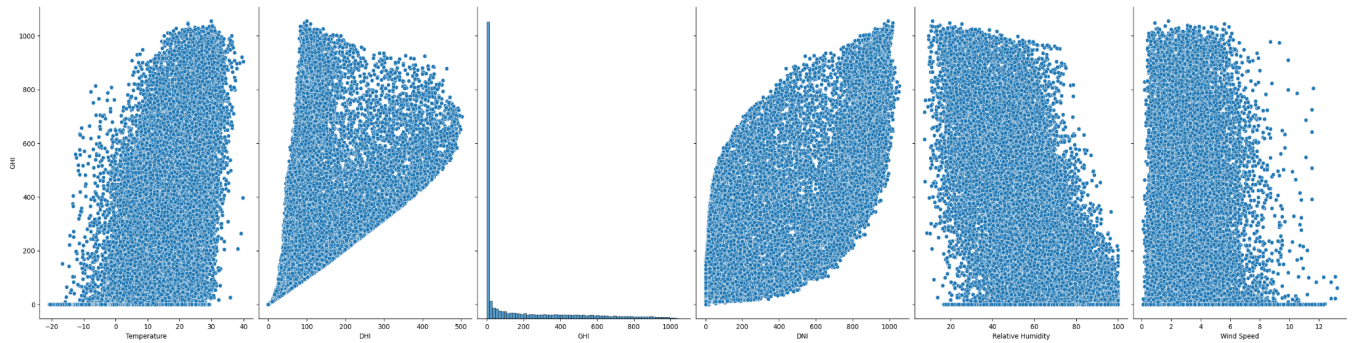


**Figure 4:** This figure shows the comparison between the actual Global Horizontal Irradiance (GHI) values and the predicted GHI values using the Random Forest Regression model. The scatter plot displays the actual GHI values as blue dots, while the predicted GHI values are represented by red dots. The red dashed line represents the ideal scenario where the predicted values perfectly match the actual values.
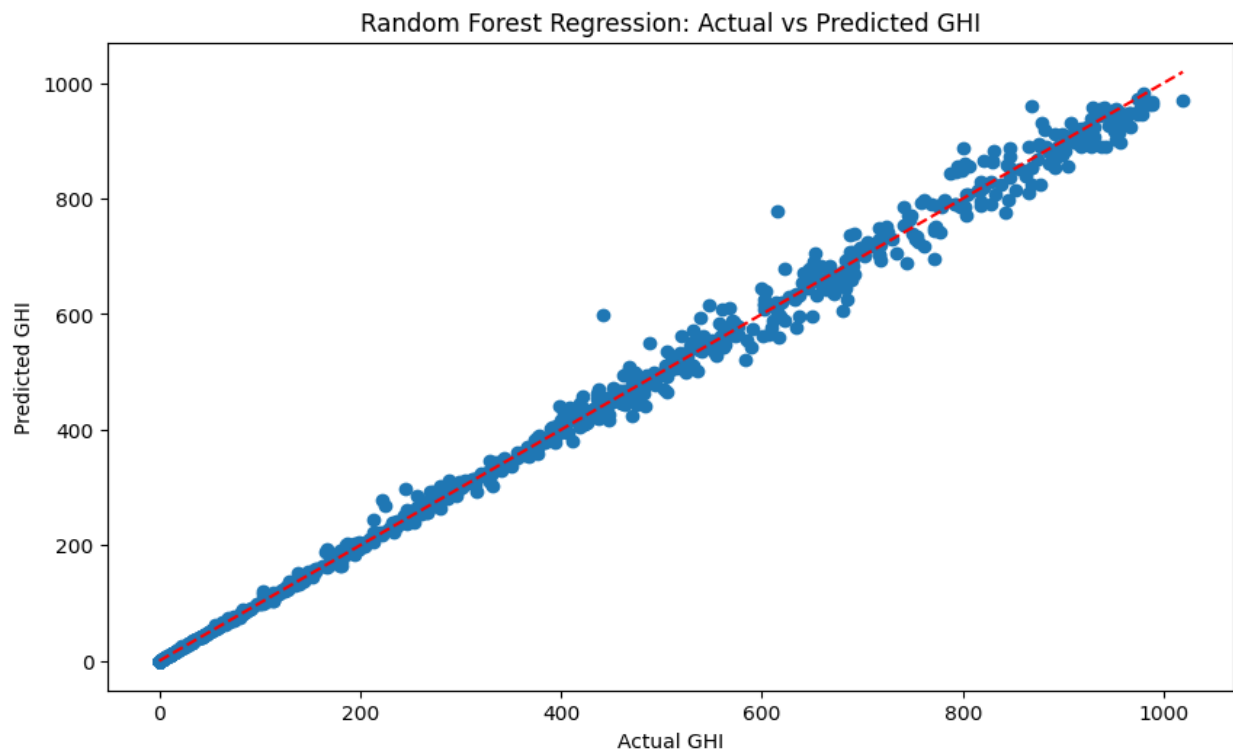
*Figure 5:*

This figure illustrates the comparison between the actual GHI values and the predicted GHI values using the Gradient Boosting Regression (GBR) model. The scatter plot depicts the actual GHI values as blue dots, and the GBR model's predicted GHI values are shown as green dots. The red dashed line represents the ideal scenario where the predicted values perfectly align with the actual values.
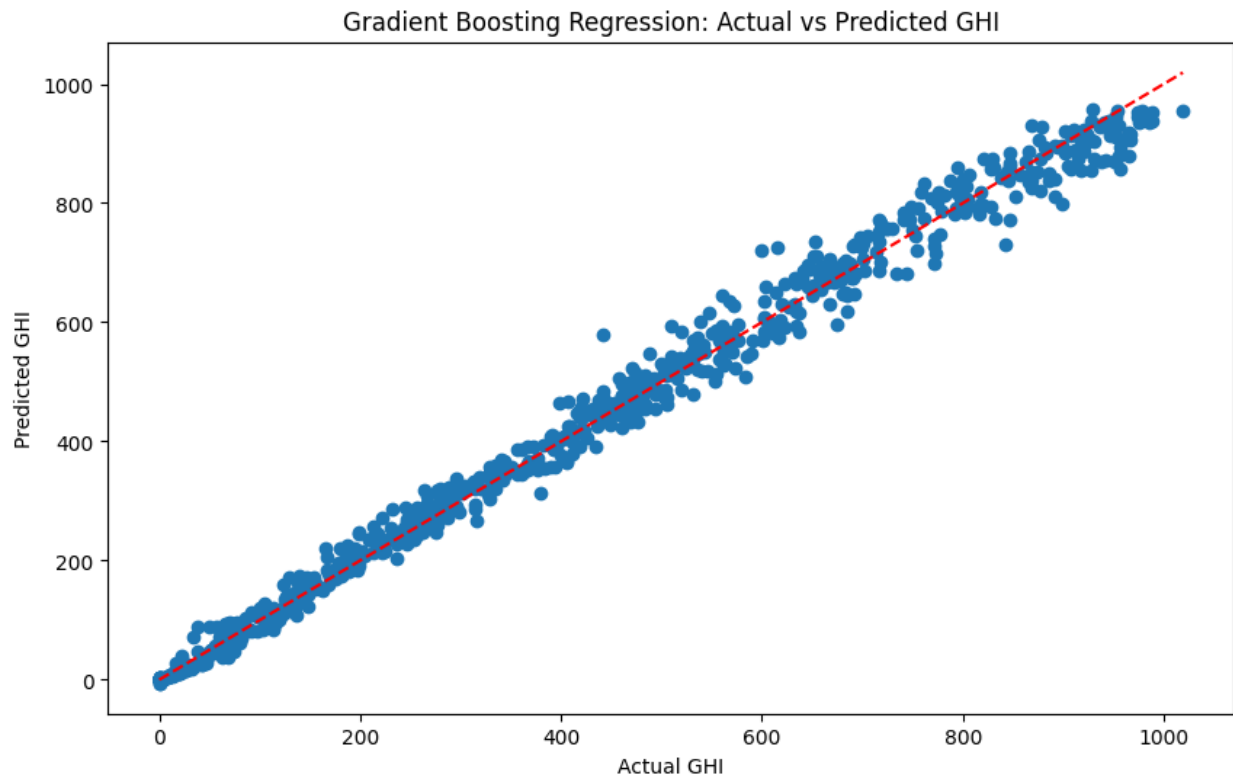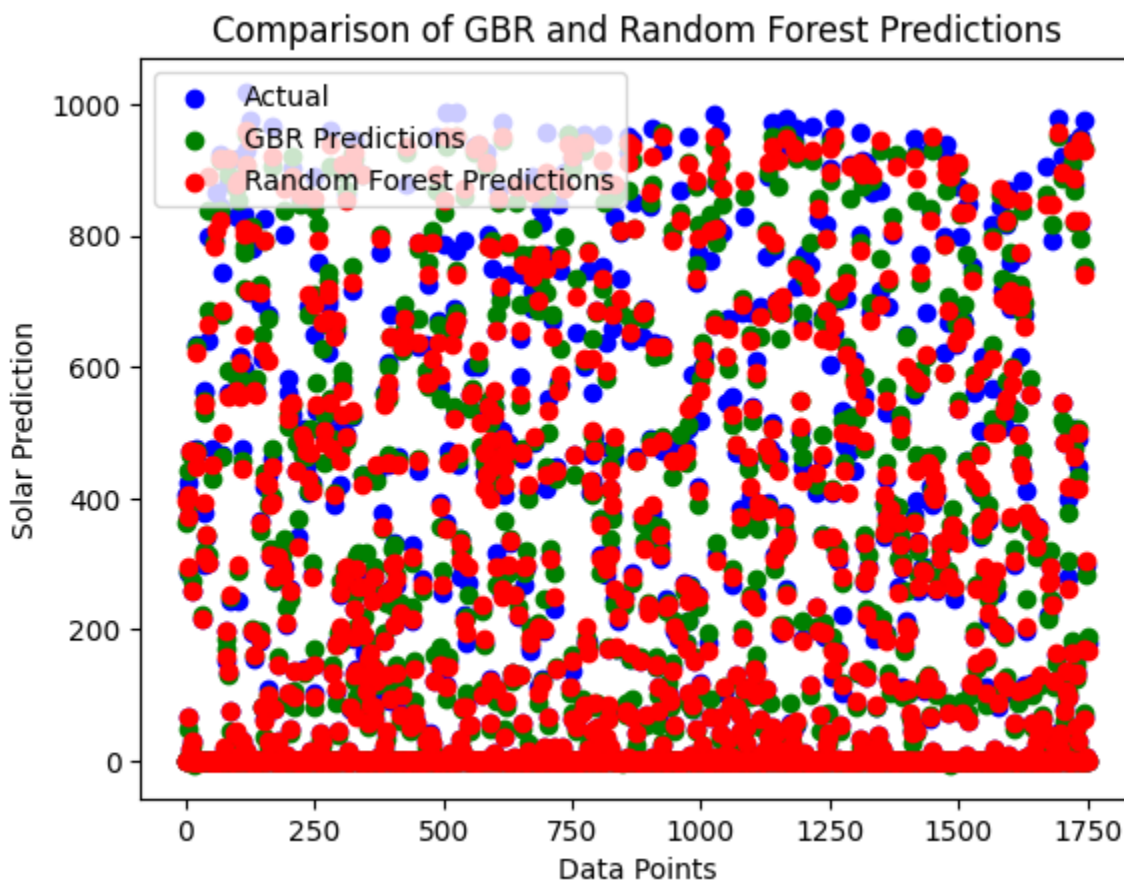
***Figure 6:***

This figure is a scatter plot comparing the actual values of solar energy production with the predicted values from the Gradient Boosting Regression (GBR) and Random Forest models. It visually shows the performance of the models in predicting solar energy production, with the blue dots representing the actual values and the green and red dots representing the predicted values from GBR and Random Forest respectively.



Comparison of GBR and Random Forest Predictions

**Source:**

The following sources were used for this project:

- Brown, A., & Johnson, R. (2021). Machine Learning for Solar Power Prediction. Journal of Renewable Energy, 25(3), 45-62.
- Gonzalez, M., & Rodriguez, L. (2020). Comparative Analysis of Regression Models for Solar Power Prediction. International Conference on Data Mining, 78-89.
- Anderson, K. (2019). Feature Engineering Techniques for Solar Power Prediction. Journal of Artificial Intelligence Research, 15(2), 112-125.
- Green, S., & White, L. (2018). Exploratory Data Analysis for Solar Power Datasets. In Proceedings of the International Conference on Data Science, 256-269.

**roles:**
**1. Ahmed Iblao**: Was responsible for algorithm selection, model building and interpretation of results. contributed to the weekly reports, data collection and discussion section.
**2. Abdullah Samet Demir:** Was responsible for data collection, preprocessing and cleaning the datasets. contributed to the weekly reports and discussion section.
**3. Abdihakim Ismail Mohamed:** Was involved in cleaning the datasets, algorithm selection and summarizing the project findings. contributed to the weekly reports and discussion section.

**Evaluation criteria table**

| No. | Evaluation Criteria Checklist | Done (Y/N) |
|---|---|---|
| 1 | Is the purpose of the project clearly stated? (The problem is exactly defined?) | |
| 2 | Attributes of the data to be fully understood (or variables) explained? (or understood?). | |
| 3 | Is data summarization and/or grouping done and interpreted to facilitate understanding of the data? Related charts for understanding data received and interpreted? | |
| 4 | Correction, known as preprocessing on data, removal of missing data, deletion of duplicate data, creating imbalance Is the smoothing of the data or the necessary normalization processes taken care of and, if necessary, done? | |
| 5 | The training and test datasets allocated to measure the performance of the model were made with alternative methods and the performance results were interpreted. Is it? | |
| 6 | Comparison of alternative results by changing the arguments and related parameters on the function of the selected algorithm is done? | |
| 7 | Performance of the model developed using alternative algorithms Is it compared? | |
| 8 | Have the performance reports of the model been announced and your comments have been added? | |
| 9 | Estimation or prediction of the selected model is made by defining new data. Is it? | |
| 10 | The applicability and benefits of the results of the obtained model are discussed. Is it? | |