



National
College of
Ireland

Configuration Manual

MSc AI for Business
Practicum 2

Cristóbal Cáceres
Student ID: 23425849

School of Computing
National College of Ireland

Supervisor:
Dr Muslim Jameel Syed

National College of Ireland

MSc AI for Business

Project Submission Sheet

School of Computing



National College of Ireland

Student Name: Cristóbal Cáceres Ortúzar

Student ID: 23425849

Programme: Ai For Business Year: 2024 - 2025

Module: Practicum 2

Lecturer: Dr Muslim Jameel Syed

Submission Due Date: 11th August 2024

Project Title: Optimizing Green Hydrogen Production: An AI-Based Approach

Word Count: xxx Page Count: 15

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature: [Handwritten Signature]

Date: 11th August 2024

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Table with 2 columns: Instruction and checkbox. Row 1: Attach a completed copy of this sheet to each project (including multiple copies) [] Row 2: Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies). [] Row 3: You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. []

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Cristóbal Cáceres Ortúzar

Student ID: 23425849

1. System Requirements

RAM: 16GB

OS: Windows 10

Processor: intel evo i7 g

Technology required: Python

The project needs Python 3.10+ along with Jupyter Notebook or Google Colab for execution.

Main Libraries: pandas, numpy, matplotlib, seaborn, scikit-learn, skopt, shap, joblib, statsmodels

2. Essential Steps to Follow

Step 1: Data Loading

The renewable_hydrogen_dataset_2535.csv file was accessed by Pandas. The file contains 2,535 entries where multiple independent variables like solar and wind power along with electrolyzer features and environmental parameters match with the target variable Hydrogen_Production_kg/day. The data gets loaded into a DataFrame which provides high-speed operations for modification and statistical analysis before ML integration. (Urhan et al., 2025; Naveena et al., 2024).

Step 2: Data Preprocessing

The model required numerical data so all non-numerical entries were discarded. The independent variables (X) received separate treatment from the dependent variable (y) that measured hydrogen production on a daily basis. The correlation matrix displayed variable relationships, but VIF calculations helped detect multicollinearity in the dataset. The model training requires independent non-redundant features. (Mukelabai, Barbour & Blanchard, 2024; Kabir et al., 2023).

Step 3: Data Splitting

The Scikit-learn library enabled us to separate our dataset into training and testing portions at 80% and 20% using `train_test_split()`. The training data received model fitting from the models which used the testing data to validate generalization results. The `random_state` parameter was fixed at 42 to guarantee that the results will be replicable. (Schröer et al., 2021).

Step 4: Model Comparison

A comparison of four regression models took place between Linear Regression and Support Vector Regression (SVR) with RBF kernel and Gradient Boosting Regressor and Random Forest Regressor. The models received their training data from the training set before moving on to test their performance with the test set. The assessment of model performance included Root Mean Squared Error (RMSE) together with Mean Absolute Error (MAE) and R² Score for quantitative model comparison. (Mallala et al., 2025; Motiramani et al., 2025).

Step 5: Hyperparameter Optimization

The process of model tuning involved two distinct stages. The Random Forest model received its first set of optimal parameters through Grid Search with 5-fold cross-validation. BayesSearchCV performed a second optimization step that balanced the exploration and exploitation of hyperparameters. The method produces faster computations than exhaustive search while selecting the best parameters. (Motiramani et al., 2025; Shash et al., 2025).

Step 6: Final Model Evaluation

The Random Forest model with Bayesian Optimization produced the best results which received evaluation against the test data to determine the final RMSE, MAE, and R² values. A Repeated K-Fold Cross-Validation (10 folds, 3 repeats) was performed to confirm the model's stability and generalization capabilities and generated performance statistics with standard deviation values (Mallala et al., 2025; Shanmugasundaram et al., 2025).

Step 7: Model Interpretation

The model became more understandable through the implementation of permutation importance for feature ranking. The SHAP (SHapley Additive exPlanations) values were computed to display how each feature influences prediction outcomes. The model follows XAI principles because it provides clear insights into its decision-making process. (Ahmed et al., 2024; Afzali et al., 2024).

Step 8: Sensitivity Analysis & Scenario Simulation

A sensitivity analysis was implemented to study how variations in individual features affect predicted hydrogen production. A simulation using the maximum observed values for all features was conducted to estimate potential production under optimal conditions. This step provides insights into the operational scenarios that could maximize hydrogen output (Raja et al., 2025; Wei et al., 2025).

Step 9: Error Analysis

Two diagnostic plots were generated:

Residual Plot: to check if residuals are randomly distributed around zero, indicating an unbiased model.

Actual vs. Predicted Plot – to visually compare predicted values against true outputs and assess accuracy.

The graphical tools serve to detect systematic errors while confirming model assumptions (Rai & Liu, 2025; Phan et al., 2024).

Step 10: Model Export

The final trained model received a.pkl file format from Joblib for future deployment purposes that eliminate the need for model retraining. The predicted outcomes together with residuals received.csv file format for documentation purposes and additional analysis needs. The model maintains reproducibility through this approach and enables seamless system integration (Quintanilla et al., 2025; Shahin & Simjoo, 2025).

```
In [82]: # Import Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, GridSearchCV, RepeatedKFold, cross_val_score
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.svm import SVR
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.inspection import permutation_importance
from skopt import BayesSearchCV
from joblib import dump
import shap
from statsmodels.stats.outliers_influence import variance_inflation_factor

# Load dataset from absolute path
df = pd.read_csv(r"C:\Users\Crist\OneDrive\Escritorio\FINAL_THESIS\dataset_GREEN\renewable_hydrogen_dataset_2535.csv")

# Define X and y - Keep only numeric columns
X = df.drop(columns=['Hydrogen_Production_kg/day'])
X = X.select_dtypes(include=[np.number]) # Avoid string columns
y = df['Hydrogen_Production_kg/day']

# Correlation matrix
plt.figure(figsize=(10,8))
sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm')
plt.title("Correlation Matrix")
plt.show()

# VIF calculation
X_numeric = X.select_dtypes(include=[np.number])
vif_data = pd.DataFrame()
vif_data["Feature"] = X_numeric.columns
vif_data["VIF"] = [variance_inflation_factor(X_numeric.values, i) for i in range(X_numeric.shape[1])]
print(vif_data.sort_values(by="VIF", ascending=False))
```

```

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Multi-model comparison
models = {
    'LinearRegression': LinearRegression(),
    'SVR (RBF)': SVR(kernel='rbf'),
    'GradientBoosting': GradientBoostingRegressor(random_state=42),
    'RandomForest': RandomForestRegressor(random_state=42)
}

results = []
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    mae = mean_absolute_error(y_test, y_pred)
    r2 = r2_score(y_test, y_pred)
    results.append({'Model': name, 'RMSE': rmse, 'MAE': mae, 'R2': r2})

results_df = pd.DataFrame(results).set_index('Model')
print(results_df)

# GridSearchCV for Random Forest
param_grid = {
    'n_estimators': [50, 100, 150],
    'max_depth': [3, 5, None],
    'min_samples_split': [2, 5],
    'min_samples_leaf': [1, 2]
}

grid_search = GridSearchCV(RandomForestRegressor(random_state=42), param_grid, cv=5, scoring='r2', n_jobs=-1, verbose=2)
grid_search.fit(X_train, y_train)
best_grid_model = grid_search.best_estimator_
print("Best GridSearch Params:", grid_search.best_params_)

```

```

# Bayesian Optimization
param_space = {
    'n_estimators': (50, 300),
    'max_depth': (3, 20),
    'min_samples_split': (2, 10),
    'min_samples_leaf': (1, 5)
}

bayes_search = BayesSearchCV(
    RandomForestRegressor(random_state=42),
    param_space,
    n_iter=50,
    cv=5,
    scoring='r2',
    random_state=42,
    n_jobs=-1,
    verbose=2
)

bayes_search.fit(X_train, y_train)
best_bayes_model = bayes_search.best_estimator_
print("Best Bayesian Params:", bayes_search.best_params_)

# Final evaluation
y_pred_final = best_bayes_model.predict(X_test)
rmse = np.sqrt(mean_squared_error(y_test, y_pred_final))
mae = mean_absolute_error(y_test, y_pred_final)
r2 = r2_score(y_test, y_pred_final)
print(f"Final Model - RMSE: {rmse:.3f}, MAE: {mae:.3f}, R2: {r2:.3f}")

# Repeated 10-Fold CV
rkf = RepeatedKFold(n_splits=10, n_repeats=3, random_state=42)
scores = cross_val_score(best_bayes_model, X, y, scoring='r2', cv=rkf, n_jobs=-1)
print(f"Repeated 10-Fold CV Mean R2: {scores.mean():.3f}, Std: {scores.std():.3f}")

```

```

# Permutation Importance
perm_importance = permutation_importance(best_bayes_model, X_test, y_test, n_repeats=10, random_state=42, n_jobs=-1)
sorted_idx = perm_importance.importances_mean.argsort()
plt.figure(figsize=(8,6))
plt.barh(X_test.columns[sorted_idx], perm_importance.importances_mean[sorted_idx])
plt.xlabel("Permutation Importance")
plt.title("Feature Importance")
plt.show()

# SHAP analysis
explainer = shap.Explainer(best_bayes_model, X_train)
shap_values = explainer(X_test)
shap.summary_plot(shap_values, X_test)

# Sensitivity Analysis
def sensitivity_analysis(variable, model, X_test, step=50):
    var_range = np.linspace(X_test[variable].min(), X_test[variable].max(), step)
    mean_preds = []
    for val in var_range:
        temp_X = X_test.copy()
        temp_X[variable] = val
        pred = model.predict(temp_X).mean()
        mean_preds.append(pred)
    plt.figure(figsize=(8,5))
    plt.plot(var_range, mean_preds)
    plt.xlabel(variable)
    plt.ylabel('Predicted Hydrogen Production (kg/day)')
    plt.title(f'Sensitivity Analysis: {variable}')
    plt.grid(True)
    plt.show()

for var in X_test.columns:
    sensitivity_analysis(var, best_bayes_model, X_test)

# Scenario Simulation (Fix for warning)
max_values = X.max()
pred_max_scenario = best_bayes_model.predict(pd.DataFrame([max_values], columns=X.columns))[0]
print(f"Predicted Production for Optimal Scenario: {pred_max_scenario:.2f} kg/day")

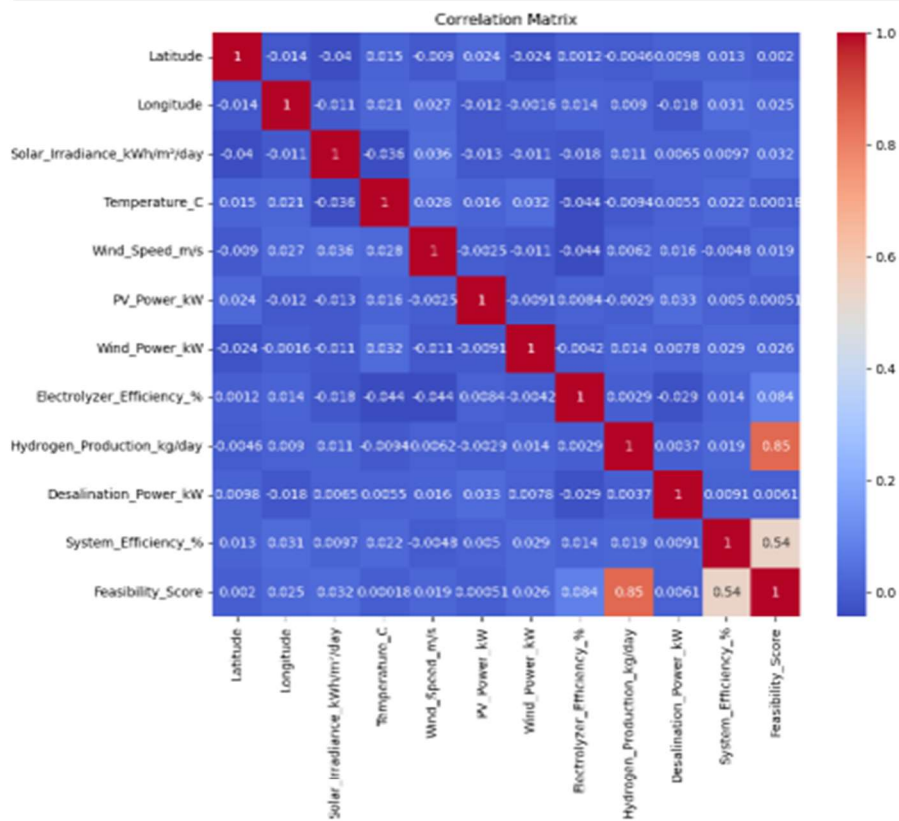
# Residual Plot
residuals = y_test - y_pred_final
plt.figure(figsize=(8,5))
sns.scatterplot(x=y_pred_final, y=residuals)
plt.axhline(y=0, color='r', linestyle='--')
plt.xlabel('Predicted')
plt.ylabel('Residuals')
plt.title('Residual Plot')
plt.show()

# Actual vs Predicted
plt.figure(figsize=(8,5))
sns.scatterplot(x=y_test, y=y_pred_final)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--')
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.title('Actual vs Predicted')
plt.show()

# Export model and predictions
dump(best_bayes_model, 'final_random_forest_model.pkl')
final_predictions_df = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred_final, 'Residuals': residuals})
final_predictions_df.to_csv('final_predictions_results.csv', index=False)
print("✅ Model and results exported.")

```

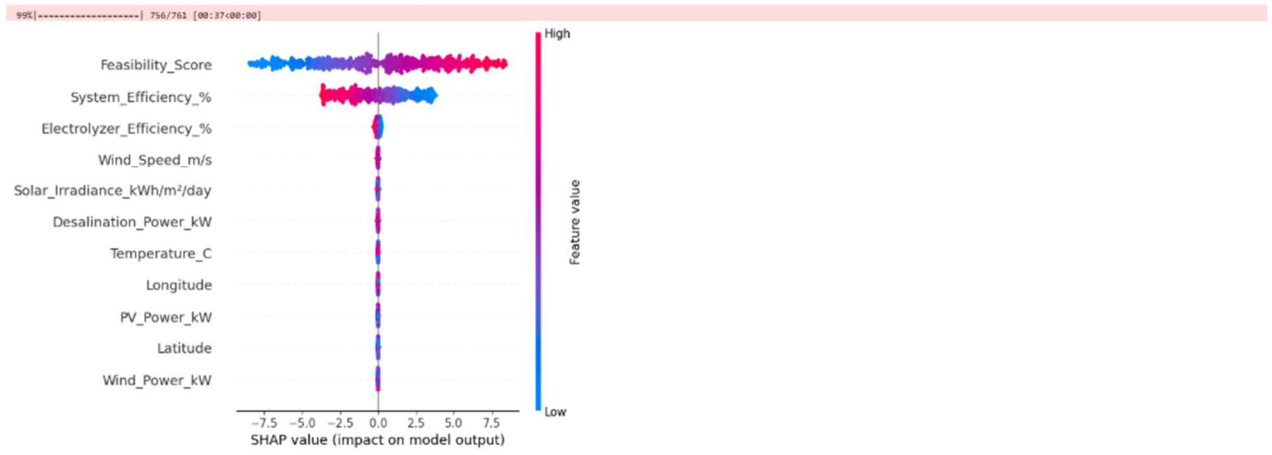
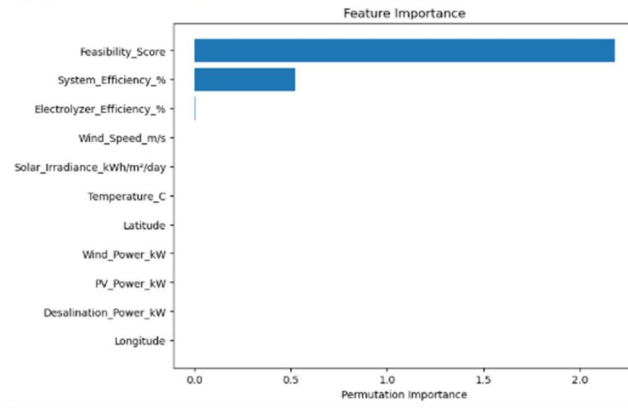
Output:

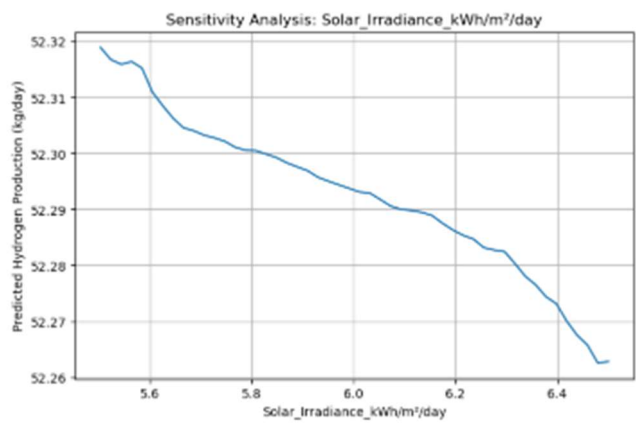
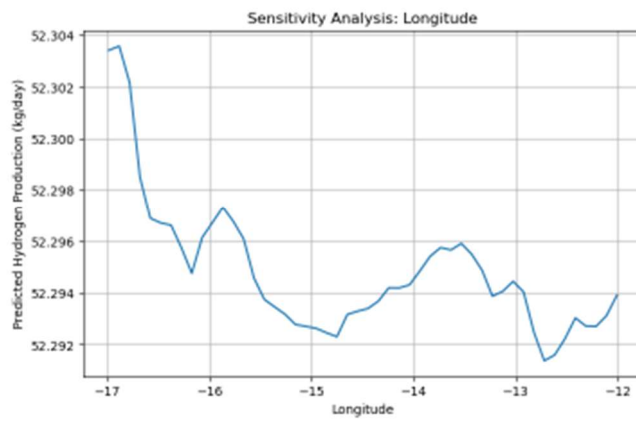
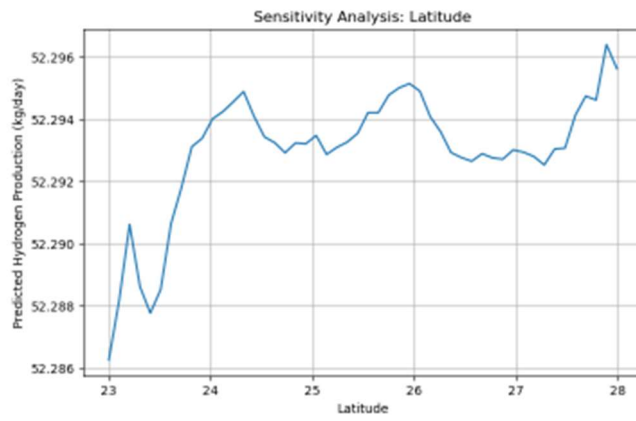


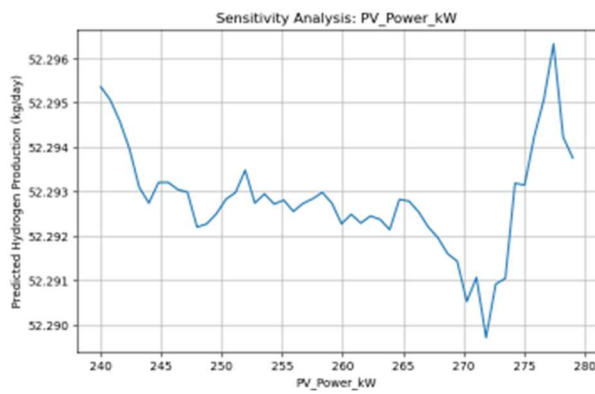
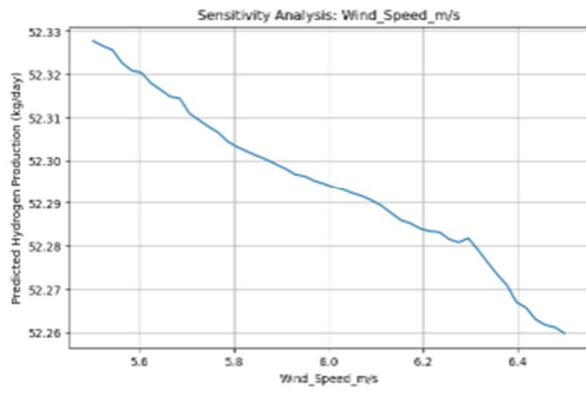
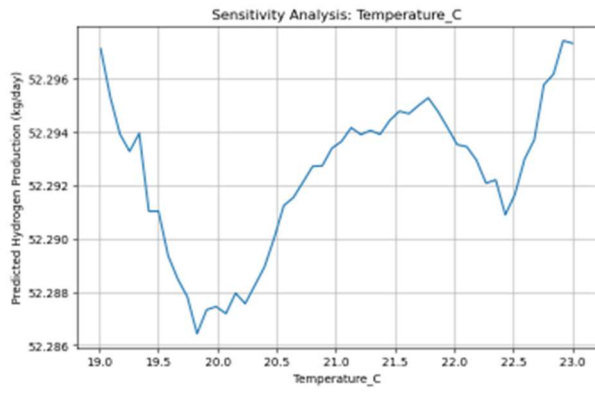
	Feature	VIF
7	Electrolyzer_Efficiency_%	2188.884217
9	System_Efficiency_%	2811.783346
10	Feasibility_Score	1820.912889
6	Wind_Power_kW	685.883762
5	PV_Power_kW	493.895847
2	Solar_Irradiance_kWh/m ² /day	418.973676
4	Wind_Speed_m/s	414.110623
3	Temperature_C	318.771386
8	Latitude	306.388912
8	Desalination_Power_kW	126.898358
1	Longitude	180.953628

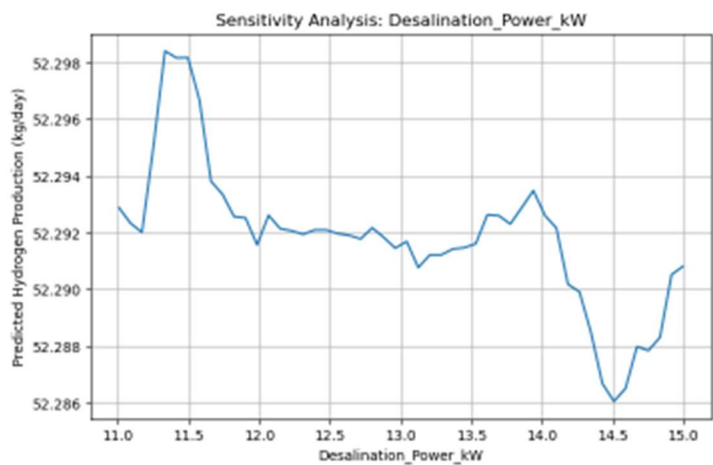
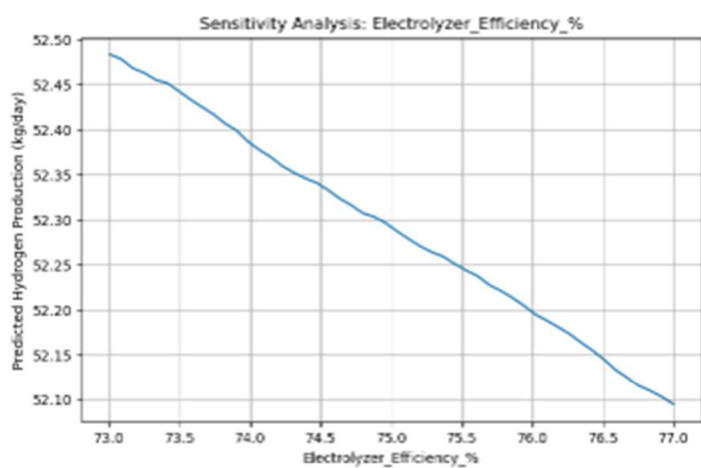
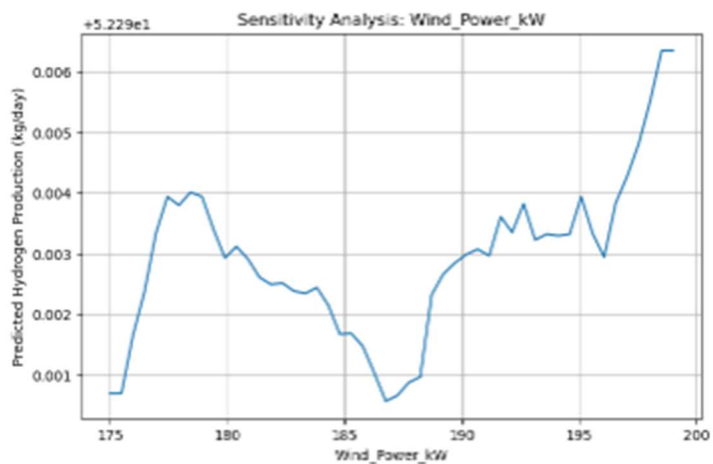
R2

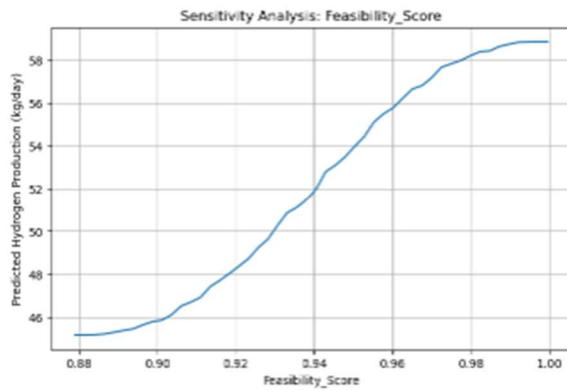
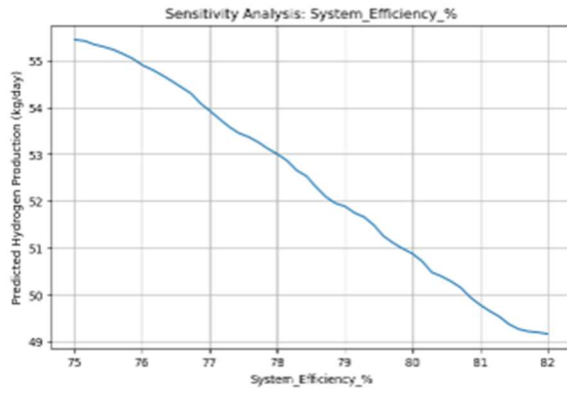
Repeated 10-Fold CV Mean R²: 0.993, Std: 0.001



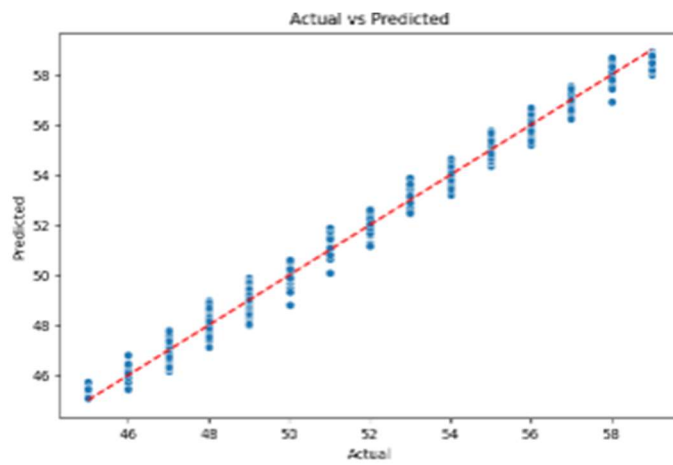
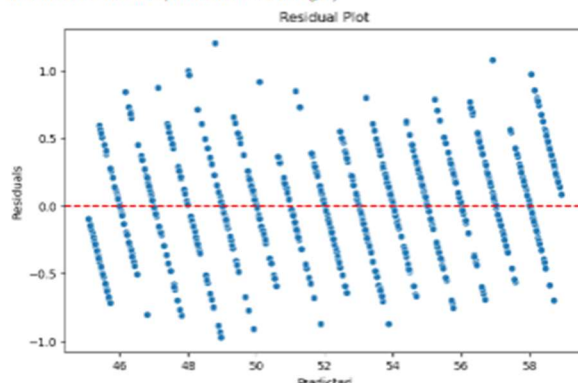








Predicted Production for Optimal Scenario: 58.25 kg/day



Model and results exported.

REFERENCES.

1. Urhan, B.B. et al., 2025. Predicting green hydrogen production using electrolyzers driven by photovoltaic panels and wind turbines based on machine learning techniques: A pathway to on-site hydrogen refuelling stations. *International Journal of Hydrogen Energy*, 101, pp.1421–1438. Available at: <https://doi.org/10.1016/j.ijhydene.2025.01.017>
 2. Naveena, K. et al., 2024. Elevating sustainability with a multi-renewable hydrogen generation system empowered by machine learning and multi-objective optimization. *Measurement: Sensors*, 33, 101192. Available at: <https://doi.org/10.1016/j.measen.2024.101192>
 3. Mukelabai, M.D., Barbour, E.R. & Blanchard, R.E., 2024. Modeling and optimization of renewable hydrogen systems: A systematic methodological review and machine learning integration. *Energy and AI*, 18, 100455. Available at: <https://doi.org/10.1016/j.egyai.2024.100455>
 4. Kabir, M.M. et al., 2023. Machine learning-based prediction and optimization of green hydrogen production technologies from water industries for a circular economy. *Desalination*, 567, 116992. Available at: <https://doi.org/10.1016/j.desal.2023.116992>
 5. Schröer, C. et al., 2021. A systematic literature review on applying CRISP-DM process model. *Procedia Computer Science*, 181, pp.526–534. Available at: <https://doi.org/10.1016/j.procs.2021.01.199>
 6. Mallala, B. et al., 2025. Forecasting global sustainable energy from renewable sources using random forest algorithm. *Results in Engineering*, 25, p.103789. Available at: <https://doi.org/10.1016/j.rineng.2024.103789>
 7. Motiramani, M. et al., 2025. AI-ML techniques for green hydrogen: A comprehensive review. *Next Energy*, 8, 100252. Available at: <https://doi.org/10.1016/j.nxener.2025.100252>
 8. Shash, A.Y. et al., 2025. Computational methods, Artificial Intelligence, Modeling, and Simulation applications in green hydrogen production through water electrolysis: A review. *Hydrogen*, 6(2), 21. Available at: <https://doi.org/10.3390/hydrogen6020021>
 9. Shanmugasundaram, S. et al., 2025. A review on green hydrogen production pathways and optimization techniques. *Process Safety and Environmental Protection*, 197, 107070. Available at: <https://doi.org/10.1016/j.psep.2025.107070>
 10. Ahmed, R. et al., 2024. An explainable AI for green hydrogen production: A deep learning regression model. *International Journal of Hydrogen Energy*, 83, pp.1226–1242. Available at: <https://doi.org/10.1016/j.ijhydene.2024.08.064>
 11. Wei, J. et al., 2025. Artificial intelligence applications in hydrogen system: Advancing renewable energy utilization for global hydrogen economy and sustainability goals. *International Journal of Hydrogen Energy*, 122, pp.359–373. Available at: <https://doi.org/10.1016/j.ijhydene.2025.03.350>
 12. Raja, I.B. et al., 2025. Regional variability in the performance of Solar-Green Hydrogen Hybrid Energy Systems (SGHHES): Synergistic enviro-economic analysis and evaluation across six climatic zones using multi-criteria decision analysis. *International Journal of Hydrogen Energy*, 138, pp.681–693. Available at: <https://doi.org/10.1016/j.ijhydene.2025.05.193>
 13. Rai, A. & Liu, J., 2025. Optimizing solar-electrolysis for green hydrogen production: A novel spatiotemporal attention framework (STAF) for solar-electrolysis prediction and economic viability analysis. *International Journal of Hydrogen Energy*, 142, pp.723–738. Available at: <https://doi.org/10.1016/j.ijhydene.2025.03.385>
 14. Phan, T.P. et al., 2024. Dynamic-based artificial intelligence model for simulation and optimization of the single chamber anode brush microbial electrolysis cell. *Fuel*, 375, 132629. Available at: <https://doi.org/10.1016/j.fuel.2024.132629>
 15. Quintanilla, P. et al., 2025. Artificial intelligence and robotics in the hydrogen lifecycle: A systematic review. *International Journal of Hydrogen Energy*, 113, pp.801–817. Available at: <https://doi.org/10.1016/j.ijhydene.2025.03.016>
 16. Shahin, M. & Simjoo, M., 2025. Potential applications of innovative AI-based tools in hydrogen energy development: Leveraging large language model technologies. *International Journal of Hydrogen Energy*, 102, pp.918–936. Available at: <https://doi.org/10.1016/j.ijhydene.2025.01.066>
- DATA SET: https://www.kaggle.com/datasets/ziya07/renewable-hydrogen-production-dataset/data?select=renewable_hydrogen_dataset_2535.csv