# Solar Power Generation Forecasting

MALSAWMZUALA KHIANGTE
IIT BHUBANESWAR

# Problem Statement:

The irregular nature of solar energy generation creates big challenges for keeping the power grid stable and managing energy efficiently. Sudden changes in generation can cause energy waste, higher operating costs, and greater dependence on less eco-friendly backup power sources.

This project builds a machine learning model tailored to a solar park in Mizoram to accurately predict its power output. This helps operators plan better, manage resources efficiently, and integrate solar energy smoothly into the grid — reducing energy waste and supporting Mizoram's shift toward clean energy.

# Project Workflow

**DATA COLLECTION
&
DATA CLEANING**

**FEATURE
ENGINEERING**

**MODEL TRAINING
&
EVALUATION**

**DEPLOYMENT**

# Data Acquisition

*Raw Data*



*Cleaning data*

```python
time_pattern = re.compile(r'(\d{2}:\d{2})')        # matches times like 06:15
date_pattern = re.compile(r'Date\s*:\s*(\d{2}/\d{2}/\d{4})', re.IGNORECASE)
number_pattern = re.compile(r'-?\d+\.\d+|-?\d+')

for idx, row in df.iterrows():
    # join row into a single string for searching
    row_str = ' '.join(row.dropna().astype(str).values
    # detect date line like "Date : 01/04/2023"
    date_match = date_pattern.search(row_str)
    if date_match:
        current_date = date_match.group(1)
        continue
    if current_date is None:
        # skip lines before the first date
        continue
    # try to extract the first time (prefer 'From' time)
    time_match = time_pattern.search(row_str)
    if not time_match:
        # no time on this row — skip
        continue
    time_str = time_match.group(1)  # e.g. '06:15'
    # Extract numeric tokens from the row string
    nums = [float(x) for x in number_pattern.findall(row_str)]
    # Heuristic to find Total MW
    mw_value = None
    if len(nums) >= 2:
        mw_value = nums[-2]
    if mw_value is None or np.isnan(mw_value):
        continue
    # compute energy for this 15-min block (MWh)
    energy_block_mwh = mw_value * 0.25
    # compose full timestamp
    try:
        ts = pd.to_datetime(f"{current_date} {time_str}", dayfirst=True, format="%d/%m/%Y %H:%
    except Exception:
        ts = pd.to_datetime(f"{current_date} {time_str}", dayfirst=True)

    records.append({
        'timestamp': ts,
        'Total_MW': mw_value,
        'Block_Energy_MWh': energy_block_mwh
    })
```

*Cleaning data*

```python
    # create DataFrame of blocks
    blocks_df = pd.DataFrame(records)
    blocks_df = blocks_df.set_index('timestamp').sort_index()

    # resample to hourly
    hourly = blocks_df.resample('h').agg({
        'Block_Energy_MWh': 'sum',      # total energy produced in that hour (MWh)
        'Total_MW': ['mean', 'max']     # average MW and max MW during that hour
    })

    # flatten MultiIndex columns
    hourly.columns = ['MWh', 'Avg_MW', 'Max_MW']
    hourly = hourly.reset_index()

    # Save individual monthly results
    #hourly.to_csv(output_filename, index=False)
    print(f"Successfully processed and saved '{output_filename}'")

    return hourly

# --- Main script execution ---

# List of months to process
months_to_process = [
    'April', 'May', 'June', 'July', 'August',
    'September', 'October', 'November', 'December'
    ]

all_dataframes = []
for month in months_to_process:
    monthly_df = process_energy_file(month)
    if not monthly_df.empty:
        all_dataframes.append(monthly_df)

# Combine all the monthly data into a single DataFrame
if all_dataframes:
    final_combined_df = pd.concat(all_dataframes, ignore_index=True)

    print("\n--- All months combined successfully! ---")
    final_combined_df.to_csv('all_months_hourly_energy.csv', index=False)
    print("Final combined data saved to 'all_months_hourly_energy.csv'")
    print("\nFinal DataFrame Info:")
    print(final_combined_df.info())
else:
    print("\nNo data was processed. The final DataFrame is empty.")
```

# Data Cleaning

- Aggregated 9 months (Apr-Dec 2023) of raw 15-minute solar production logs into a single dataset.

- Standardized park's IST logs and weather's UTC data to a common timezone.

- Upscaled 15-minute production data to hourly averages to match the weather data's granularity.

- Merged the two clean, time-synced datasets (Production + Weather) on the new hourly timestamp.

- Complete Sanity Check. Checked for physically impossible values (e.g., humidity > 100, negative ghi).

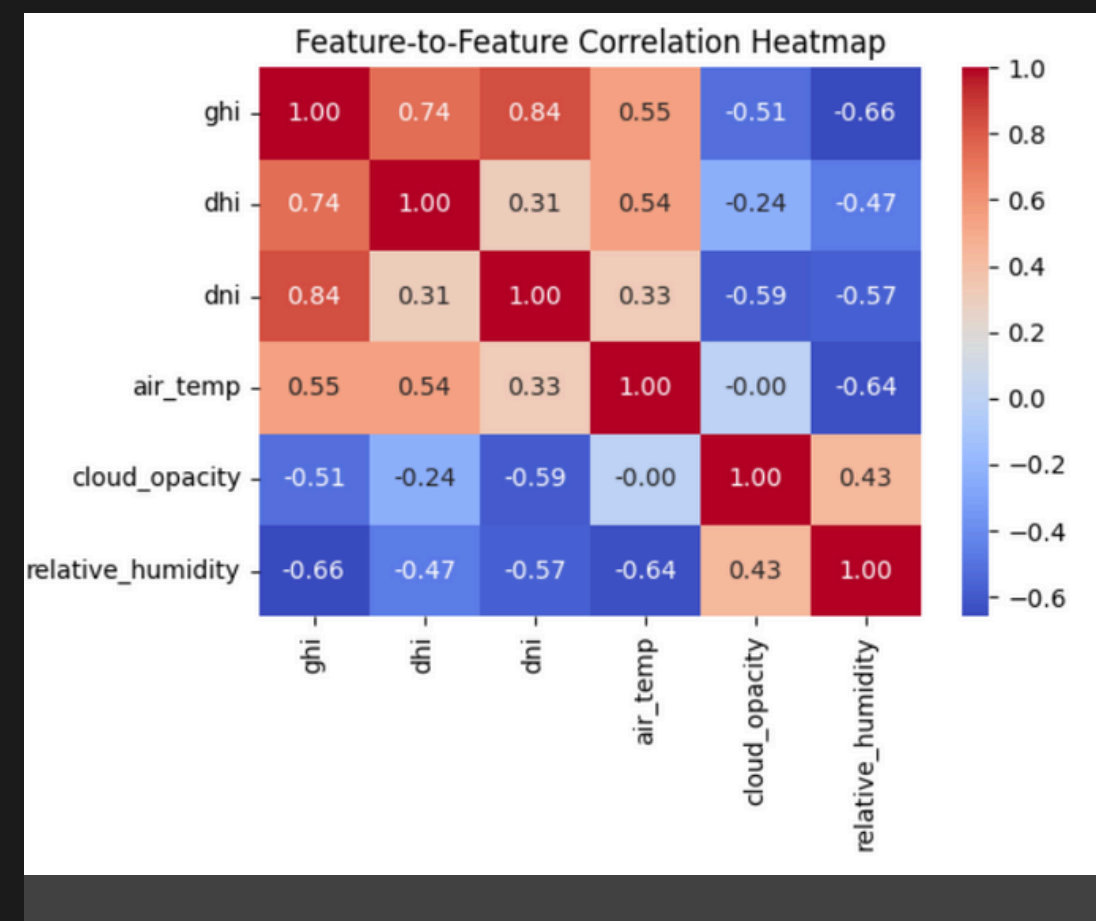- Confirmed all feature data (.describe() table) was within realistic physical bounds.



*Final dataset*

# Feature Engineering



**TARGET CORRELATION ANALYSIS.**

**MULTICOLLINEARITY ANALYSIS**

**SCATTERPLOT OF FINAL FEATURES WITH TARGET VARIABLE**

# Feature Engineering



Box Plot of MWh to Identify Outliers

| COUNT | 6558.000000 |
|---|---|
| MEAN | 4.130625 |
| STD | 5.883984 |
| MIN | 0.000000 |
| 25% | 0.000000 |
| 50% | 0.082150 |
| 75% | 7.743325 |
| MAX | 22.500000 |

*THE OUTLIER BOUNDARY IS: 19.3583 MWH*
*THERE ARE 6 OUTLIERS OUT OF 6558 TOTAL ROWS.*

◇ **Finding**: Statistical analysis (1.5 * IQR rule) identified 6 data points as "upper outliers" (values > 19.36 MWh).

◇ **Investigation:** These high values were cross-referenced with their corresponding weather data.

◇ **Conclusion:** These outliers are NOT errors or bad data. They are real, valid, peak-production events—representing the sunniest, clearest, and most productive hours in the dataset.

◇ **Decision: KEEP ALL 6 OUTLIERS.**

◇ **Reasons:** The model must learn what conditions lead to peak output. Removing these points would prevent the model from ever predicting a high-production day.

# Model Training And Evaluation

- First, the feature matrix (X) was created by dropping the target variable (MWh) and all non-numeric date/time columns. The target vector (y) was then set as the MWh column.

- This dataset was then split into an 80% training set and a 20% test set to prepare for modeling.

- StandardScaler was fit only on the training dataset to avoid data leakage and to normalize the feature scales before model training

- Finally, a range of models (Linear Regression, Decision Tree, etc.) were trained on this scaled data and evaluated using $R^2$ score, MAE, and RMSE to measure their performance.

- A comparative evaluation of all models was conducted. The Random Forest, once optimized with GridSearchCV, was identified as the champion model, yielding the highest predictive accuracy by achieving the top $R^2$ score and the lowest MAE.
  This final, optimized model was then saved as a .pkl file for future use in live predictions.

```python
X = merged.drop(['MWh', 'date', 'time', 'month'], axis=1)
y=merged['MWh']

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test= train_test_split(X, y, test_size=0.2,random_state=42)

from sklearn.preprocessing import StandardScaler
scaler=StandardScaler()

X_train_scaled=scaler.fit_transform(X_train)
X_test_scaled=scaler.transform(X_test)
from sklearn.linear_model import LinearRegression
lr_model=LinearRegression()
lr_model.fit(X_train_scaled,y_train)
y_pred_lr=lr_model.predict(X_test_scaled)


from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error

r2=r2_score(y_test,y_pred_lr)
mae=mean_absolute_error(y_test,y_pred_lr)
mse=mean_squared_error(y_test,y_pred_lr)
rmse=np.sqrt(mse)
```

```python
from sklearn.ensemble import RandomForestRegressor
rf_model=RandomForestRegressor(random_state=4|2)
rf_model.fit(X_train_scaled,y_train)
y_pred_rf=rf_model.predict(X_test_scaled)

from sklearn.model_selection import GridSearchCV
param_grid = {
    'n_estimators': [100, 200],          # How many trees in the forest
    'max_depth': [None, 10, 20],         # How deep each tree can go
    'min_samples_leaf': [1, 2, 4]        # Min samples required at a leaf node
}
grid_search=GridSearchCV(estimator=RandomForestRegressor(random_state=42),param_grid=param_grid,
                         cv=10,n_jobs=-1,scoring='neg_mean_squared_error',verbose=1)
grid_search.fit(X_train_scaled,y_train)
print(f'Best Params found : {grid_search.best_params_}')
```

```
Best Params found : {'max_depth': 10, 'min_samples_leaf': 4, 'n_estimators': 100}
```

```python
best_rf_model=grid_search.best_estimator_
y_pred_best_rf=best_rf_model.predict(X_test_scaled)

r2=r2_score(y_test,y_pred_best_rf)
mae=mean_absolute_error(y_test,y_pred_best_rf)
mse=mean_squared_error(y_test,y_pred_best_rf)
rmse=np.sqrt(mse)
print("\n--- Tuned Random Forest Performance ---")
print(f'R2 Score : {r2:.4f}')
print(f"Mean Absolute Error : {mae:.4f}")
print(f"Root Mean Squared Error : {rmse:.4f} MWh")
```

```
--- Tuned Random Forest Performance ---
R2 Score : 0.9072
Mean Absolute Error : 0.8808
Root Mean Squared Error : 1.7879 MWh
```

# Model Training And Evaluation



Comparison of Regression Model Performance

# Model deployment

- **Load Tools:** The system first loads the three essential "production-ready" files:
    1. The saved solar_prediction_model.pkl
    2. The saved scaler.pkl
    3. The model_columns.json (to ensure column order)

- **Fetch Live Data:** A script calls the Solcast API in real-time to get the next 24 hours of weather forecast data for the Mizoram park's exact coordinates.

- **Prepare Features:** The raw JSON data from the API is processed:
    1. Timestamps are converted from UTC to the local timezone (IST).
    2. Data is formatted to exactly match the feature columns the model was trained on.

- **Scale & Predict:**
    1. The loaded scaler is used to .transform() the new weather data.
    2. The scaled data is fed into the model.predict() function to get the MWh forecasts.

- Clean & Display: The result is a clean, 24-hour forecast table, which can be saved to a CSV

| | Local_Time | Predicted_MWh |
|---|---|---|
| 0 | 2025-10-30 16:00 | 1.8527 |
| 1 | 2025-10-30 17:00 | 0.0155 |
| 2 | 2025-10-30 18:00 | 0.0107 |
| 3 | 2025-10-30 19:00 | 0.0101 |
| 4 | 2025-10-30 20:00 | 0.0093 |
| 5 | 2025-10-30 21:00 | 0.0361 |
| 6 | 2025-10-30 22:00 | 0.0289 |
| 7 | 2025-10-30 23:00 | 0.0040 |
| 8 | 2025-10-31 00:00 | 0.0182 |
| 9 | 2025-10-31 01:00 | 0.0042 |
| 10 | 2025-10-31 02:00 | 0.0044 |
| 11 | 2025-10-31 03:00 | 0.0044 |
| 12 | 2025-10-31 04:00 | 0.0044 |
| 13 | 2025-10-31 05:00 | 0.0044 |
| 14 | 2025-10-31 06:00 | 2.5741 |
| 15 | 2025-10-31 07:00 | 5.4134 |
| 16 | 2025-10-31 08:00 | 10.6889 |
| 17 | 2025-10-31 09:00 | 13.3560 |
| 18 | 2025-10-31 10:00 | 14.3939 |
| 19 | 2025-10-31 11:00 | 13.9482 |
| 20 | 2025-10-31 12:00 | 12.6150 |
| 21 | 2025-10-31 13:00 | 10.5143 |
| 22 | 2025-10-31 14:00 | 7.9699 |

# Conclusion

In conclusion, this project built a successful forecasting tool for the Mizoram P&ED that turns unpredictable solar power into a reliable asset. This model directly saves money by reducing costly emergency power purchases. The Tuned Random Forest is a strong baseline, and the future roadmap is clear.

The immediate next step is to re-validate all models using a TimeSeriesSplit to get a true, robust performance score that prevents any data leakage. Following that, engineering cyclical (sin/cos) features for time of day will further enhance the model's accuracy, delivering even greater financial savings.