# Climate Connect

## Solar Power Generation from the plant.

**Objectives:** To Predict the Power Generation from the Forecast Data.

**Introduction:** Design a profitable solar power plant by analyzing the radiation and metrological data. You will benefit in the operation by reduced maintenance costs and high efficiency.

The sun radiation on earth surface combines **Direct normal irradiation (DNI)** and **Diffuse Horizontal Irradiation (DHI)** .Both are linked in the formula for **Global Horizontal Irradiation (GHI).**

$$GHI + DHI + DNI\text{-}\cos(\theta)$$

**GTI** is total amount of direct and diffuse irradiation received from above by a tilted surface.

From the Given Data **GHI** in known to you and how much Power is generated with help of GHI.

From the Research  the Parameter that influence the Solar Energy Production **Air-temperature, wind-speed, wind-direction, humidity etc.** details which contains in Weather Actuals file.

## Approach:

### Feature Engineering

1. Converting variables into proper **date time** format.

```
In [11]: result['datetime_utc'] = pd.to_datetime(result['datetime_utc'])
         result['datetime_local'] = pd.to_datetime(result['datetime_local'])
         result['sunrise'] = pd.to_datetime(result['sunrise'])
         result['sunset'] = pd.to_datetime(result['sunset'])
         result['updated_at'] = pd.to_datetime(result['updated_at'])
```

2. Identifying **Null values** and **Drop** the variables.

```
In [12]: result.isnull().sum()
         humidity                0
         dew_point               0
         wind_bearing            0
         wind_speed              0
         wind_chill           9875
         wind_gust               0
         heat_index           9875
         pressure                0
         qpf                  9875
         uv_index                0
         snow                 9875
         pop                  9875
         fctcode              9875
         ozone                   0
         precip_accumulation  9875
         precip_intensity        0
         precip_probability      0
         precip_type          7129
         visibility              0
         sunrise                 0
```

```
In [13]: result.drop(['plant_id','precip_type','precip_accumulation','fctcode','pop',\
                 'snow','qpf','heat_index','wind_chill','ghi','gti'],axis=1,inplace=True)
```

3. **Heat map** describing the correlation between the variables and drop the variables highly correlated with each other.
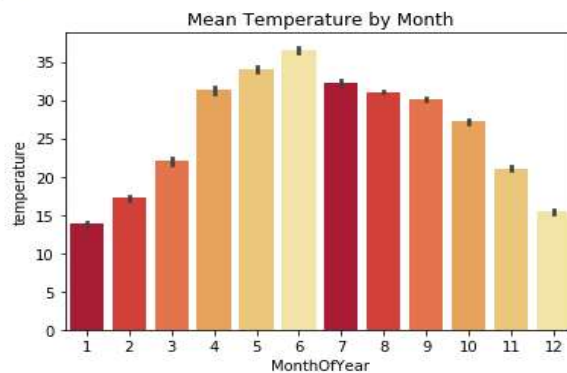


4. Calculating the time duration between **Sunrise** and **Sunset**.

```
In [17]: result['MonthOfYear'] = result['datetime_utc'].dt.month
         result['DayOfYear'] = result['datetime_utc'].dt.dayofyear
         result['WeekOfYear'] = result['datetime_utc'].dt.weekofyear
         result['TimeOfDay(h)'] = result['datetime_utc'].dt.hour
         result['TimeOfDay(m)'] = result['datetime_utc'].dt.hour*60 + result['datetime_utc'].dt.minute
         result['TimeOfDay(s)'] = result['datetime_utc'].dt.hour*60*60 + result['datetime_utc'].dt.minute*60
```

```
In [18]: result['DayLength(s)'] = result['sunset'].dt.hour*60*60 + result['sunset'].dt.minute*60 + result['sunset'].dt.second \
                                 - result['sunrise'].dt.hour*60*60 - result['sunrise'].dt.minute*60 - result['sunrise'].dt.second
```

5. visualizing the **temperature** vary in month. It describes in **month of 4, 5, 6(summer)** keeps on increasing and decreases **after summer over**.

```
In [41]: ax2 = plt.axes()
         pal2 = sns.color_palette('YlOrRd_r')
         sns.barplot(x="MonthOfYear", y='temperature', data=result, palette=pal2, ax = ax2)
         ax2.set_title('Mean Temperature by Month')
         plt.show()
```

**Model Selection:**

1]. Select the suitable variables for prediction and split the data and applying **RandomForestRegressor (RFR)** Algorithm.

```
In [31]: X = result[['humidity','wind_speed','DayOfYear','TimeOfDay(s)','cloud_cover','DayLength(s)']]

In [32]: y = result['power']

In [33]: 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 = 0)

In [34]: from sklearn.ensemble import RandomForestRegressor
         from sklearn.model_selection import cross_val_score

In [35]: regressor = RandomForestRegressor(n_estimators = 100)
         regressor.fit(X_train, y_train)

Out[35]: RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
                               max_features='auto', max_leaf_nodes=None,
                               min_impurity_decrease=0.0, min_impurity_split=None,
                               min_samples_leaf=1, min_samples_split=2,
                               min_weight_fraction_leaf=0.0, n_estimators=100,
                               n_jobs=None, oob_score=False, random_state=None,
                               verbose=0, warm_start=False)
```

2. Applying feature Importance we get r2 score. From this we can Say **DayOfYear and DayLength(s)** are highly correlated with target variable with **r2 score** of **0.832956.**

```
for i in range(0,5):
    least_important = np.argmin(feature_importances)
    removed_columns = removed_columns.append(X_train_opt.pop(X_train_opt.columns[least_important]))
    regressor.fit(X_train_opt, y_train)
    feature_importances = regressor.feature_importances_
    accuracies = cross_val_score(estimator = regressor,
                                 X = X_train_opt,
                                 y = y_train, cv = 5,
                                 scoring = 'r2')
    r2s_opt = np.append(r2s_opt, accuracies.mean())
    models = np.append(models, ", ".join(list(X_train_opt)))

feature_selection = pd.DataFrame({'Features':models,'r2 Score':r2s_opt})
feature_selection.head()
```

Out[36]:

| | Features | r2 Score |
|---|---|---|
| 0 | wind_speed, DayOfYear, TimeOfDay(s), cloud_cov... | 0.693153 |
| 1 | wind_speed, DayOfYear, TimeOfDay(s), DayLength(s) | 0.706911 |
| 2 | wind_speed, DayOfYear, DayLength(s) | 0.707771 |
| 3 | DayOfYear, DayLength(s) | 0.832956 |
| 4 | DayLength(s) | 0.646000 |

3. Finally the Best variables which predict the power supply are **DayOfLength(s), DayOfYear, and TimeOfDay** with **r2 score** of **0.9234** even it performs better on test data with r2 score more than train data with r2 score of **0.92577** on test data as well.

```
In [45]: X = result[['DayLength(s)','DayOfYear','TimeOfDay(s)']]
         y=result['power']
         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
         regressor = RandomForestRegressor(n_estimators = 100)
         regressor.fit(X_train, y_train)
         accuracies = cross_val_score(estimator = regressor, X = X_train,y = y_train, cv = 10, scoring = 'r2')
         accuracy = accuracies.mean()
         print('r2 = {}'.format(accuracy))

         r2 = 0.9234834872606632

In [42]: from sklearn.metrics import r2_score
         y_pred = regressor.predict(X_test)
         r_squared = r2_score(y_test, y_pred)
         print('r2 = {}'.format(r_squared))

         r2 = 0.9257704125646082
```

**Conclusion:**

From the Given Data DayLength in seconds and DayOfYear and TimeOfDay in seconds gives best predictive model for prediction of Solar Power Generation.