## Problem Statement:

The Global Power Plant Database is a comprehensive, open source database of power plants around the world. It centralizes power plant data to make it easier to navigate, compare and draw insights for one's own analysis. The database covers approximately 35,000 power plants from 167 countries and includes thermal plants (e.g. coal, gas, oil, nuclear, biomass, waste, geothermal) and renewables (e.g. hydro, wind, solar). Each power plant is geolocated and entries contain information on plant capacity, generation, ownership, and fuel type. It will be continuously updated as data becomes available.

## Importing Required Library

```
In [1]: import numpy as np
        import pandas as pd
        import matplotlib.pyplot as plt
        import seaborn as sns
        import pickle
        from scipy.stats import zscore
        import scikitplot as skplt
        pd.set_option('display.max_columns', None) # # For display maximum column
        from sklearn.preprocessing import StandardScaler, LabelEncoder
        from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_sco
        from sklearn.tree import DecisionTreeRegressor
        from sklearn.neighbors import KNeighborsRegressor
        from sklearn.ensemble import RandomForestRegressor
        from sklearn.metrics import r2_score, mean_squared_error
        import xgboost as xgb
        %matplotlib inline

        import warnings
        warnings.filterwarnings('ignore')
```
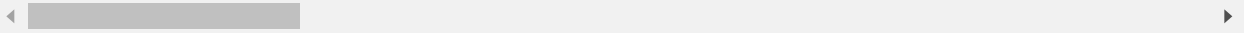
## Reading Data

```
In [2]: df = pd.read_csv(r"C:\Users\Kushal Arya\Desktop\csv file\database_IND.csv")
        df.head()
```

Out[2]:

| | country | country_long | name | gppd_idnr | capacity_mw | latitude | longitude | primary_fuel |
|---|---------|--------------|------|-----------|-------------|----------|-----------|--------------|
| 0 | IND | India | ACME Solar Tower | WRI1020239 | 2.5 | 28.1839 | 73.2407 | Solar |
| 1 | IND | India | ADITYA CEMENT WORKS | WRI1019881 | 98.0 | 24.7663 | 74.6090 | Coal |
| 2 | IND | India | AES Saurashtra Windfarms | WRI1026669 | 39.2 | 21.9038 | 69.3732 | Wind |
| 3 | IND | India | AGARTALA GT | IND0000001 | 135.0 | 23.8712 | 91.3602 | Gas |
| 4 | IND | India | AKALTARA TPP | IND0000002 | 1800.0 | 21.9603 | 82.4091 | Coal |

## Check no of row and column

```
In [3]: print('No of Rows and Columns ----->', df.shape )
```

No of Rows and Columns -----> (908, 25)

## Checking for Null values

```
In [4]: print('----------------------\n')
        print(df.isnull().sum())
        print('\n----------------------')
```

```
----------------------

country                        0
country_long                   0
name                           0
gppd_idnr                      0
capacity_mw                    0
latitude                      46
longitude                     46
primary_fuel                   0
other_fuel1                  709
other_fuel2                  907
other_fuel3                  908
commissioning_year           380
owner                        566
source                         0
url                            0
geolocation_source            19
wepp_id                      908
year_of_capacity_data        388
generation_gwh_2013          524
generation_gwh_2014          507
generation_gwh_2015          483
generation_gwh_2016          471
generation_gwh_2017          465
generation_data_source       458
estimated_generation_gwh     908
dtype: int64

----------------------
```

**There is null value**

## Information about dataset

```
In [5]: print('-------------------------------------------------------\n')
        print(df.info())
        print('-------------------------------------------------------')
```

```
        -------------------------------------------------------

        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 908 entries, 0 to 907
        Data columns (total 25 columns):
         #   Column                 Non-Null Count  Dtype
        ---  ------                 --------------  -----
         0   country                908 non-null    object
         1   country_long           908 non-null    object
         2   name                   908 non-null    object
         3   gppd_idnr              908 non-null    object
         4   capacity_mw            908 non-null    float64
         5   latitude               862 non-null    float64
         6   longitude              862 non-null    float64
         7   primary_fuel           908 non-null    object
         8   other_fuel1            199 non-null    object
         9   other_fuel2            1 non-null      object
         10  other_fuel3            0 non-null      float64
         11  commissioning_year     528 non-null    float64
         12  owner                  342 non-null    object
         13  source                 908 non-null    object
         14  url                    908 non-null    object
         15  geolocation_source     889 non-null    object
         16  wepp_id                0 non-null      float64
         17  year_of_capacity_data  520 non-null    float64
         18  generation_gwh_2013    384 non-null    float64
         19  generation_gwh_2014    401 non-null    float64
         20  generation_gwh_2015    425 non-null    float64
         21  generation_gwh_2016    437 non-null    float64
         22  generation_gwh_2017    443 non-null    float64
         23  generation_data_source 450 non-null    object
         24  estimated_generation_gwh  0 non-null   float64
        dtypes: float64(13), object(12)
        memory usage: 177.5+ KB
        None
        -------------------------------------------------------
```
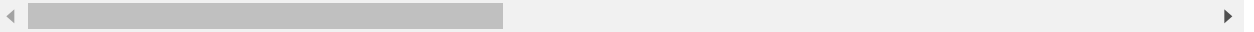
**Categorical data present in our data set**

## Statistics of Data

```
In [6]: df.describe()
```

Out[6]:

|  | capacity_mw | latitude | longitude | other_fuel3 | commissioning_year | wepp_id | year_of_ca |
|---|---|---|---|---|---|---|---|
| count | 908.000000 | 862.000000 | 862.000000 | 0.0 | 528.000000 | 0.0 | |
| mean | 321.046378 | 21.196189 | 77.447848 | NaN | 1996.876894 | NaN | |
| std | 580.221767 | 6.248627 | 4.907260 | NaN | 17.047817 | NaN | |
| min | 0.000000 | 8.168900 | 68.644700 | NaN | 1927.000000 | NaN | |
| 25% | 16.837500 | 16.771575 | 74.258975 | NaN | 1988.000000 | NaN | |
| 50% | 60.000000 | 21.778300 | 76.719250 | NaN | 2000.000000 | NaN | |
| 75% | 388.125000 | 25.516375 | 79.441475 | NaN | 2011.250000 | NaN | |
| max | 4760.000000 | 34.649000 | 95.408000 | NaN | 2018.000000 | NaN | |

**Outliers are present in our data set**

# Analysis of Null value

```
In [7]:  print('----------------------\n')
         print(df.isnull().sum())
         print('\n----------------------')
```

```
----------------------

country                        0
country_long                   0
name                           0
gppd_idnr                      0
capacity_mw                    0
latitude                      46
longitude                     46
primary_fuel                   0
other_fuel1                  709
other_fuel2                  907
other_fuel3                  908
commissioning_year           380
owner                        566
source                         0
url                            0
geolocation_source            19
wepp_id                      908
year_of_capacity_data        388
generation_gwh_2013          524
generation_gwh_2014          507
generation_gwh_2015          483
generation_gwh_2016          471
generation_gwh_2017          465
generation_data_source       458
estimated_generation_gwh     908
dtype: int64

----------------------
```

```
In [8]:  df['other_fuel1'].value_counts()
```

```
Out[8]:  Oil              196
         Gas                2
         Cogeneration       1
         Name: other_fuel1, dtype: int64
```

```
In [9]:  df['other_fuel2'].value_counts()
```

```
Out[9]:  Oil     1
         Name: other_fuel2, dtype: int64
```

```
In [10]:  df['other_fuel3'].value_counts()
```

```
Out[10]:  Series([], Name: other_fuel3, dtype: int64)
```

**Approach : We drop above other_fuel2 and other_fuel3 column because maximun null value and keep other_fuel1 and fill null value with mode**

```
In [11]: df['latitude'].value_counts()
```

```
Out[11]: 24.1917    3
         19.0004    3
         16.5697    2
         23.4639    2
         13.2450    2
                   ..
         20.9099    1
         17.2387    1
         23.5594    1
         27.3426    1
         16.5973    1
         Name: latitude, Length: 837, dtype: int64
```

```
In [12]: df['longitude'].value_counts()
```

```
Out[12]: 71.6917    4
         71.6918    3
         75.8988    3
         72.8983    3
         81.2875    3
                   ..
         80.1264    1
         76.1137    1
         74.6447    1
         86.0970    1
         79.5748    1
         Name: longitude, Length: 828, dtype: int64
```

**Approach : We fill above longitude and latitude columns with mean**

```
In [13]: df['commissioning_year'].value_counts()
```

```
Out[13]: 2013.0    28
         2015.0    26
         2012.0    23
         2016.0    21
         2014.0    17

         1958.0     1
         1949.0     1
         1954.0     1
         1956.0     1
         1927.0     1
         Name: commissioning_year, Length: 73, dtype: int64
```

**Approach: We fill commissioning_year column with mode**

```
In [14]: df['owner'].value_counts()
```

```
Out[14]: Acc Acc ltd                              4
         Sterling Agro Industries ltd.            4
         Jk Cement ltd                            4
         Tata Power Solar Systems Limited (TPREL)  3
         Ujaas Energy Limited                     3
                                                 ..
         Gm Energy ltd                            1
         Dcm & chem                               1
         REI Agro Limited                         1
         National And paper                       1
         Spr Pvt ltd                              1
         Name: owner, Length: 280, dtype: int64
```

**Approach : We fill owner column with mode**

```
In [15]: df['geolocation_source'].value_counts()
```

```
Out[15]: WRI                                    766
         Industry About                         119
         National Renewable Energy Laboratory     4
         Name: geolocation_source, dtype: int64
```

**Approach : We fill geolocation_source column with mode**

```
In [16]: df['wepp_id'].value_counts()
```

```
Out[16]: Series([], Name: wepp_id, dtype: int64)
```

**Approach : We drop wepp_id column because maximum null value**

```
In [17]: df['year_of_capacity_data'].value_counts()
```

```
Out[17]: 2018.0    520
         Name: year_of_capacity_data, dtype: int64
```

**Approach : We fill year_of_capacity_data column with bfill beacuse only one data in it**

```
In [18]:  df['generation_gwh_2013'].value_counts()

Out[18]:  0.00000        21
          14881.88000     1
          42.49645        1
          2036.00000      1
          97.73885        1
                         ..
          7229.33000      1
          657.21740       1
          507.89775       1
          8556.42400      1
          8211.00000      1
          Name: generation_gwh_2013, Length: 364, dtype: int64


In [19]:  df['generation_gwh_2014'].value_counts()

Out[19]:  0.00000        28
          6803.31250      1
          4735.13000      1
          145.81400       1
          2022.57000      1
                         ..
          6224.00000      1
          268.48085       1
          1255.73200      1
          164.32425       1
          1153.65300      1
          Name: generation_gwh_2014, Length: 374, dtype: int64


In [20]:  df['generation_gwh_2015'].value_counts()

Out[20]:  0.00000        28
          5837.76600      1
          1297.97750      1
          8076.81050      1
          1.09395         1
                         ..
          2636.86400      1
          665.19730       1
          1516.36010      1
          741.86205       1
          7130.50700      1
          Name: generation_gwh_2015, Length: 398, dtype: int64
```

```
In [21]: df['generation_gwh_2016'].value_counts()
```

```
Out[21]: 0.00000       31
         8470.57000     2
         1511.00000     2
         7.31325        1
         94.85500       1
                       ..
         433.84800      1
         283.74811      1
         259.94375      1
         403.96000      1
         307.87290      1
         Name: generation_gwh_2016, Length: 405, dtype: int64
```

```
In [22]: df['generation_gwh_2017'].value_counts()
```

```
Out[22]: 0.00000       33
         170.08530      2
         344.35955      1
         2265.47000     1
         59.43135       1
                       ..
         214.48220      1
         272.73945      1
         2887.00000     1
         12.73600       1
         158.73235      1
         Name: generation_gwh_2017, Length: 410, dtype: int64
```

**Approach : We fill above generation_gwh columns with mean**

```
In [23]: df['generation_data_source'].value_counts()
```

```
Out[23]: Central Electricity Authority    450
         Name: generation_data_source, dtype: int64
```

**Approach : We fill above generation_data_source columns with bfill because only one data in it**

```
In [24]: df['estimated_generation_gwh'].value_counts()
```

```
Out[24]: Series([], Name: estimated_generation_gwh, dtype: int64)
```

**Approach : We drop estimated_generation_gwh columns because no data in it**

## Drop Unwanted column

```
In [7]: col = ['estimated_generation_gwh', 'wepp_id', 'other_fuel2', 'other_fuel3', 'url'
```

```
In [8]: df = df.drop(col, axis = 1)
        df.head(2)
```

Out[8]:

| | country_long | name | capacity_mw | latitude | longitude | primary_fuel | other_fuel1 | commission |
|---|---|---|---|---|---|---|---|---|
| 0 | India | ACME Solar Tower | 2.5 | 28.1839 | 73.2407 | Solar | NaN | |
| 1 | India | ADITYA CEMENT WORKS | 98.0 | 24.7663 | 74.6090 | Coal | NaN | |

```
In [9]: print('After droping no of Rows and Columns  ----->', df.shape )
```

After droping no of Rows and Columns  -----> (908, 18)

## Fill NaN

```
In [10]: df = df.apply(lambda x:x.fillna(x.mean())if x.dtype == 'float' else x.fillna(x.va
```

```
In [11]: print('---------------------\n')
         print(df.isnull().sum())
         print('\n---------------------')
```

```
---------------------

country_long             0
name                     0
capacity_mw              0
latitude                 0
longitude                0
primary_fuel             0
other_fuel1              0
commissioning_year       0
owner                    0
source                   0
geolocation_source       0
year_of_capacity_data    0
generation_gwh_2013      0
generation_gwh_2014      0
generation_gwh_2015      0
generation_gwh_2016      0
generation_gwh_2017      0
generation_data_source   0
dtype: int64


---------------------
```
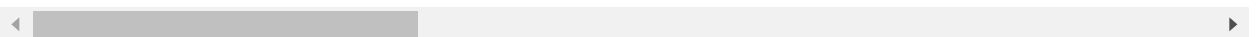
**There is no null value left**

## Analysis of data respect to Capacity MW

```
In [12]: df.head(2)
```

Out[12]:

| | country_long | name | capacity_mw | latitude | longitude | primary_fuel | other_fuel1 | commission |
|---|---|---|---|---|---|---|---|---|
| 0 | India | ACME Solar Tower | 2.5 | 28.1839 | 73.2407 | Solar | Oil | 201 |
| 1 | India | ADITYA CEMENT WORKS | 98.0 | 24.7663 | 74.6090 | Coal | Oil | 199 |

## Fule Type column
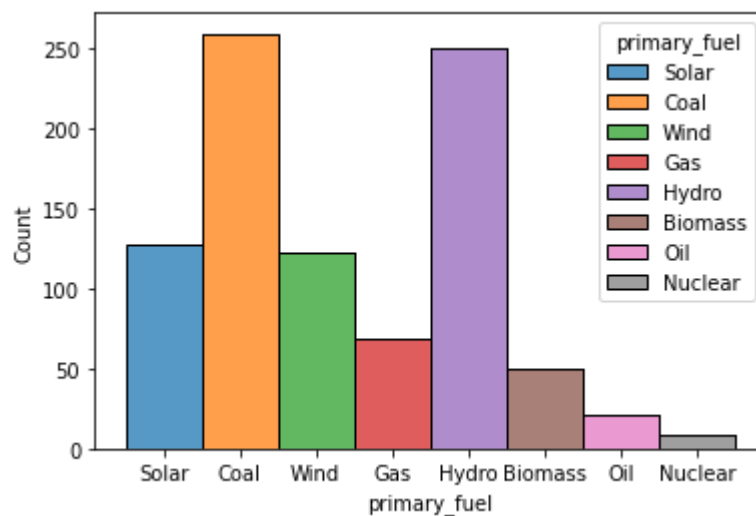
```
In [31]: df['primary_fuel'].value_counts()
```

```
Out[31]: Coal        259
         Hydro       250
         Solar       127
         Wind        123
         Gas          69
         Biomass      50
         Oil          21
         Nuclear       9
         Name: primary_fuel, dtype: int64
```

```
In [32]: sns.histplot(binwidth=0.5, x="primary_fuel", hue="primary_fuel", data=df, stat="c
         plt.show()
```



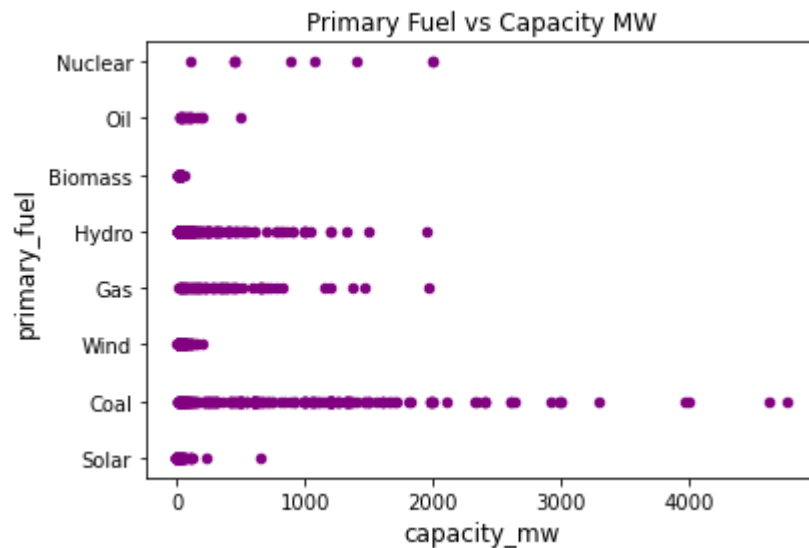**Cole type fuel used highest and Nuclear type fuel used least**

```
In [13]: df['capacity_mw'].value_counts()
```

```
Out[13]: 5.0       39
         10.0      22
         600.0     21
         15.0      20
         1200.0    19
                   ..
         26.4       1
         68.8       1
         91.8       1
         1.8        1
         816.4      1
         Name: capacity_mw, Length: 365, dtype: int64
```

```
df.plot.scatter(x = 'capacity_mw', y = 'primary_fuel', c = 'purple')
plt.xlabel('capacity_mw', fontsize = 12, c = 'black')
plt.ylabel('primary_fuel', fontsize = 12, c = 'black')
plt.title('Primary Fuel vs Capacity MW', fontsize = 12, c = 'black')
plt.show()
```



**Above plot shows Coal fuel give higest energy and Biomass give lowest energy**

## Other Fuel Type column

```
df['other_fuel1'].value_counts()
```

```
Oil             905
Gas               2
Cogeneration      1
Name: other_fuel1, dtype: int64
```

`sns.histplot(binwidth=0.5, x="other_fuel1", hue="other_fuel1", data=df, stat="cou`
`plt.show()`



**Oil fuel used most of all fuel**

## Commission Year column

In [35]: `df['commissioning_year'].value_counts()`

Out[35]:
```
1996.876894    380
2013.000000     28
2015.000000     26
2012.000000     23
2016.000000     21
              ...
1949.000000      1
1958.000000      1
1954.000000      1
1956.000000      1
1927.000000      1
Name: commissioning_year, Length: 74, dtype: int64
```

```
In [23]: df.groupby('commissioning_year')['capacity_mw'].sum().sort_values()
```

```
Out[23]: commissioning_year
         1953.0        4.00
         1937.0        5.00
         1949.0        9.30
         1956.0       10.00
         1959.0       15.00
                     ...
         2010.0    16198.00
         2012.0    16801.13
         2014.0    17468.00
         2013.0    21953.48
         2015.0    23185.50
         Name: capacity_mw, Length: 74, dtype: float64
```

```
In [20]: df.plot.scatter(x = 'capacity_mw', y= 'commissioning_year', c = 'r')
         plt.ylabel('Commissioning Year',fontsize = 14, color = 'b')
         plt.xlabel('capacity_mw',fontsize = 14, color = 'b')
         plt.title('Commissioning Year vs Capacity MW',fontsize = 14, color = 'b')
         plt.show()
```
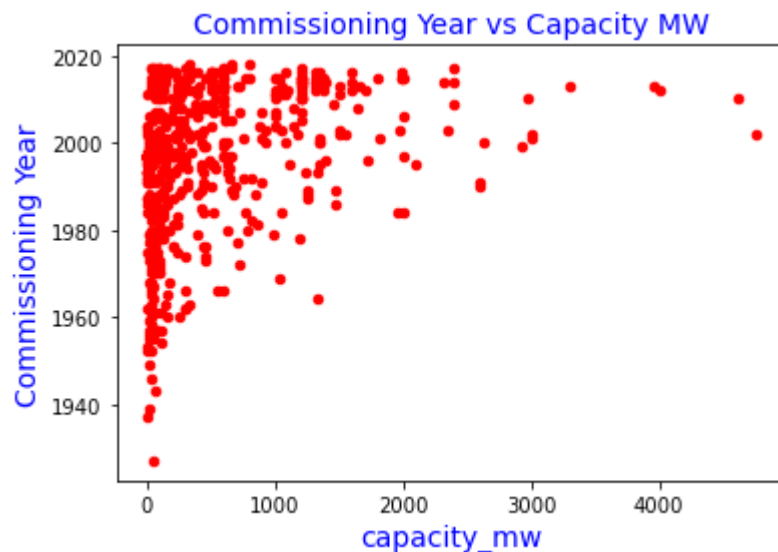


**Above plot shows early years 1927 producton of energy less than 1000 MW but Now a days it produce more than 4000 MW**

```
In [38]: df['owner'].value_counts()
```
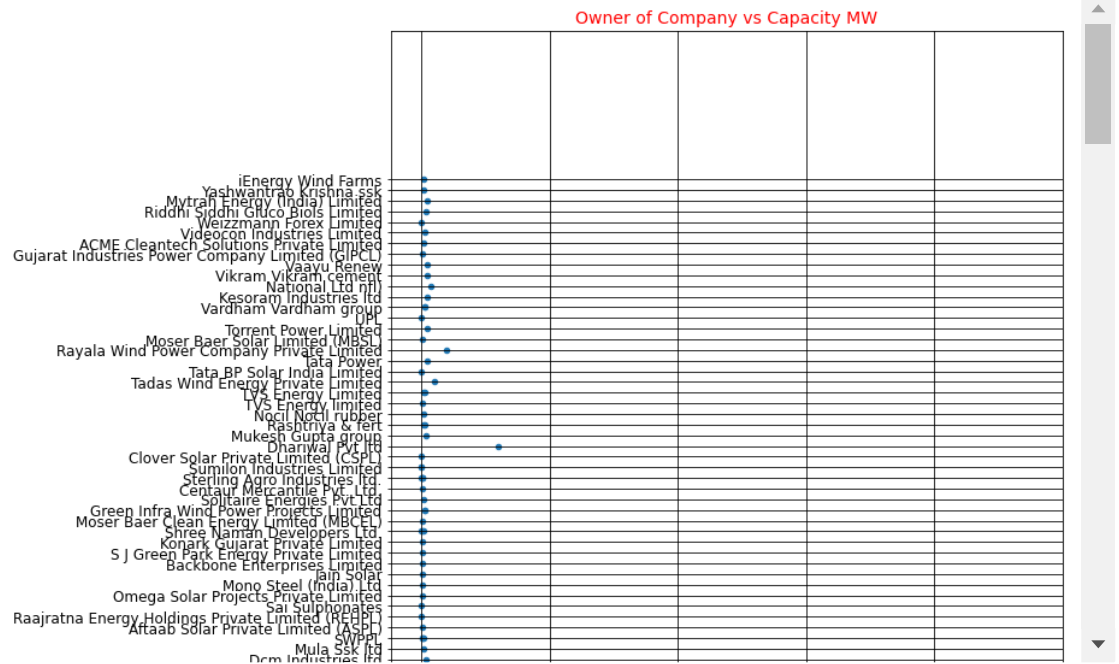
```
Out[38]: Acc Acc ltd                                570
         Sterling Agro Industries ltd.                4
         Jk Cement ltd                                4
         Tata Power Solar Systems Limited (TPREL)     3
         Ujaas Energy Limited                         3
                                                    ...
         Gm Energy ltd                                1
         Dcm & chem                                   1
         REI Agro Limited                             1
         National And paper                           1
         Spr Pvt ltd                                  1
         Name: owner, Length: 280, dtype: int64
```

```
In [39]: o = df.groupby('owner')['primary_fuel'].sum()
         o
```

```
Out[39]: owner
         ACME Cleantech Solutions Private Limited          Solar
         ACME Solar Energy                                 Solar
         AES                                                Wind
         AEW Infratech Private Limited                     Solar
         Abellon CleanEnergy Limited                       Solar
                                                            ...
         West Coast Paper Mills Ltd.                         Gas
         Yashwantrao Krishna ssk                         Biomass
         Ym Ssk ltd                                      Biomass
         Zamil New Delhi Infrastructure Private Limited    Solar
         iEnergy Wind Farms                                 Wind
         Name: primary_fuel, Length: 280, dtype: object
```

```
In [25]: df.plot.scatter(y = 'owner', x = 'capacity_mw',figsize = (10,50), rot = 360, font
         plt.grid(c = 'black')
         plt.ylabel('Owner',fontsize = 14, color = 'r')
         plt.xlabel('capacity_mw',fontsize = 14, color = 'r')
         plt.title('Owner of Company vs Capacity MW',fontsize = 14, color = 'r')
         plt.show()
```



**Above plot show Jk Cement company consumption energy is more than 4000 MW and other company consume energy less than 1000 MW**

## Co-ordinates

```
In [26]: import geopandas as gpd
         from shapely.geometry import Point, Polygon
         import descartes
         from geopandas import GeoDataFrame
         from pyproj import CRS
```

```
In [27]: geometry = [Point(xy) for xy in zip(df['longitude'], df['latitude'])]
         gdf = GeoDataFrame(df, geometry=geometry)

         #this is a simple map that goes with geopandas
         world = gpd.read_file(gpd.datasets.get_path('naturalearth_lowres'))
         gdf.plot(ax=world.plot(figsize=(10, 6)), marker='o', color='red', markersize=15);
```



**Above plot shows all cordinate belong to indian**

```
In [28]:  indmap = gpd.read_file(r"C:\Users\Kushal Arya\Desktop\csv file\map.shx")
          fig,ax = plt.subplots(figsize = (15,15))
          indmap.plot(ax = ax)
```

Out[28]:  <AxesSubplot:>

```
In [29]:  crs=CRS('EPSG:4326').to_proj4()
```

```
In [30]:  geometry = [Point(xy) for xy in zip(df['longitude'], df['latitude'])]
          geometry[:3]
```

```
Out[30]:  [<shapely.geometry.point.Point at 0x200d8d86fa0>,
           <shapely.geometry.point.Point at 0x200d94afb50>,
           <shapely.geometry.point.Point at 0x200d94af370>]
```

```
In [31]:  geo_df = gpd.GeoDataFrame(df, crs = crs, geometry = geometry)
          geo_df.head()
```

Out[31]:

| | country_long | name | capacity_mw | latitude | longitude | primary_fuel | other_fuel1 | commissic |
|---|---|---|---|---|---|---|---|---|
| 0 | India | ACME Solar Tower | 2.5 | 28.1839 | 73.2407 | Solar | Oil | 20 |
| 1 | India | ADITYA CEMENT WORKS | 98.0 | 24.7663 | 74.6090 | Coal | Oil | 19 |
| 2 | India | AES Saurashtra Windfarms | 39.2 | 21.9038 | 69.3732 | Wind | Oil | 19 |
| 3 | India | AGARTALA GT | 135.0 | 23.8712 | 91.3602 | Gas | Oil | 20 |
| 4 | India | AKALTARA TPP | 1800.0 | 21.9603 | 82.4091 | Coal | Oil | 20 |

```python
try:
    fig, ax = plt.subplots(figsize = (15,15))
    indmap.plot(ax = ax, alpha = 0.4, color = 'grey')
    geo_df[geo_df['capacity_mw']].plot(ax = ax, markersize = 20, color = 'b', mar
    geo_df[geo_df['capacity_mw']].plot(ax = ax, markersize = 20, color = '^', mar
    plt.legend(prop = {'size': 15})
except:
    pass
```

```python
In [34]: import plotly
         import plotly.graph_objects as go
         import chart_studio.plotly as py
         from plotly.offline import iplot
```

```
In [35]: data = dict(type = 'choropleth',
               locations = df['geometry'],
               locationmode = 'country names',
               z = df['capacity_mw'],
               text = df['geometry'],
               colorscale = 'YlGnBu',
               autocolorscale = False,
             )
layout = dict(geo = {'scope':'asia'})

diagram = go.Figure(data = [data], layout = layout)
iplot(diagram)
```



```
In [52]: from geopy.geocoders import Nominatim
```

```
In [53]: # initialize Nominatim API
geolocator = Nominatim(user_agent="geoapiExercises")
```

```
In [54]: # Latitude & Longitude input
         Latitude = "28.1839"
         Longitude = "73.2407"

         location = geolocator.reverse(Latitude+","+Longitude)

         address = location.raw['address']

         # traverse the data
         city = address.get('city', '')
         state = address.get('state', '')
         country = address.get('country', '')
         code = address.get('country_code')
         zipcode = address.get('postcode')
         print('City : ', city)
         print('State : ', state)
         print('Country : ', country)
         print('Zip Code : ', zipcode)
```

```
City :
State :  Rajasthan
Country :  India
Zip Code :  None
```

## Drop Columns

```
In [36]: col = ['geometry', 'generation_data_source', 'source', 'owner', 'name', 'country_
```

```
In [37]: df = df.drop(col, axis = 1)
         df.head(2)
```

Out[37]:

| | capacity_mw | latitude | longitude | primary_fuel | other_fuel1 | commissioning_year | year_of_capacit |
|---|---|---|---|---|---|---|---|
| 0 | 2.5 | 28.1839 | 73.2407 | Solar | Oil | 2011.000000 | |
| 1 | 98.0 | 24.7663 | 74.6090 | Coal | Oil | 1996.876894 | |

```
In [38]: print('After droping no of Rows and Columns  ----->', df.shape )
```

```
After droping no of Rows and Columns  -----> (908, 12)
```

## Encoding columns

```
In [39]: le = LabelEncoder()
```

```
In [40]: df['primary_fuel'] = le.fit_transform(df['primary_fuel'])
```

```
In [41]: df['other_fuel1'] = le.fit_transform(df['other_fuel1'])
```

```
In [42]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 908 entries, 0 to 907
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   capacity_mw           908 non-null    float64
 1   latitude              908 non-null    float64
 2   longitude             908 non-null    float64
 3   primary_fuel          908 non-null    int32
 4   other_fuel1           908 non-null    int32
 5   commissioning_year    908 non-null    float64
 6   year_of_capacity_data 908 non-null    float64
 7   generation_gwh_2013   908 non-null    float64
 8   generation_gwh_2014   908 non-null    float64
 9   generation_gwh_2015   908 non-null    float64
 10  generation_gwh_2016   908 non-null    float64
 11  generation_gwh_2017   908 non-null    float64
dtypes: float64(10), int32(2)
memory usage: 78.2 KB
```

```
In [62]: df.head(2)
```

Out[62]:

| | capacity_mw | latitude | longitude | primary_fuel | other_fuel1 | commissioning_year | year_of_capacit |
|---|---|---|---|---|---|---|---|
| 0 | 2.5 | 28.1839 | 73.2407 | 6 | 2 | 2011.000000 | |
| 1 | 98.0 | 24.7663 | 74.6090 | 1 | 2 | 1996.876894 | |

**All columns are encoded**

## Data distribution and checking outliers and skewness

```python
In [43]: print('---------------------')
         print('Distribution Plot :- ')
         print('---------------------')

         plt.figure(figsize = (20,25))
         plotnumber = 1

         for column in df:
             if plotnumber <=12:
                 ax = plt.subplot(4,3, plotnumber)
                 sns.distplot(df[column])
                 plt.xlabel(column, fontsize = 20)
             plotnumber +=1
         plt.tight_layout()
```

```
---------------------
Distribution Plot :-
---------------------
```

```
In [44]: df.skew()
```

Out[44]:
```
capacity_mw              3.193257
latitude                -0.147391
longitude                1.129836
primary_fuel             0.471141
other_fuel1            -20.464435
commissioning_year      -1.383330
year_of_capacity_data    0.000000
generation_gwh_2013      5.241491
generation_gwh_2014      5.041961
generation_gwh_2015      5.367370
generation_gwh_2016      5.071758
generation_gwh_2017      5.111938
dtype: float64
```

**Data has outliers and skewness**

## Corelation of Feature vs Label using Heat map

```
In [65]: print('-----------')
         print('Heat Map :-')
         print('-----------')
         df_corr = df.corr().abs()

         plt.figure(figsize = (22,16))
         sns.heatmap(df_corr, vmin = -1, annot = True, square = True, center = 0, fmt = '.
         plt.tight_layout()
```

```
-----------
Heat Map :-
-----------
```

**generation_gwh_2013,generation_gwh_2014 higest relation**

## Removing Outliers using Zscore

In [45]:
```python
# with std 3 Lets see the stats

z_score = zscore(df[['longitude', 'generation_gwh_2013', 'generation_gwh_2014',
abs_z_score = np.abs(z_score)

filtering_entry = (abs_z_score < 3).all(axis = 1)

df = df[filtering_entry]

df.describe()
```

Out[45]:

| | capacity_mw | latitude | longitude | primary_fuel | other_fuel1 | commissioning_year | year_c |
|---|---|---|---|---|---|---|---|
| count | 871.000000 | 871.000000 | 871.000000 | 871.000000 | 871.000000 | 871.000000 | |
| mean | 272.583767 | 21.089931 | 77.036419 | 3.257176 | 1.995408 | 1996.855732 | |
| std | 439.479660 | 6.145269 | 4.197722 | 2.303519 | 0.082918 | 13.017081 | |
| min | 0.000000 | 8.168900 | 68.644700 | 0.000000 | 0.000000 | 1927.000000 | |
| 25% | 16.500000 | 16.899050 | 74.329400 | 1.000000 | 2.000000 | 1996.876894 | |
| 50% | 50.700000 | 21.196189 | 76.760600 | 3.000000 | 2.000000 | 1996.876894 | |
| 75% | 330.000000 | 25.114850 | 78.922500 | 6.000000 | 2.000000 | 2003.000000 | |
| max | 2400.000000 | 34.649000 | 91.565000 | 7.000000 | 2.000000 | 2018.000000 | |

## Checking Outliers and skewness removed or not

```
In [46]:  # Let' see outliers are removed in columns or not.
          print('----------------------')
          print('Distribution Plot :- ')
          print('----------------------')

          plt.figure(figsize = (20,25))
          plotnumber = 1

          for column in df:
              if plotnumber <=12:
                  ax = plt.subplot(4,3, plotnumber)
                  sns.distplot(df[column])
                  plt.xlabel(column, fontsize = 20)
              plotnumber +=1
          plt.tight_layout()
```

```
----------------------
Distribution Plot :-
----------------------
```

**Outliers are removed**

# Spliting Dataset into features and labels

```
In [47]: x = df.drop('capacity_mw', axis = 1)
         y = df. capacity_mw
         print('Data has been splited')
```

Data has been splited

```
In [49]:  # Let' see relation between features and labels.
          print('------------------------')
          print('Scatter Plot :-')
          print('------------------------')

          plt.figure(figsize = (20,25), facecolor = 'yellow')
          plotnumber = 1
          for column in x:
              if plotnumber <=12:
                  ax = plt.subplot(4,3, plotnumber)
                  plt.scatter(x[column],y, c = 'g')
                  plt.xlabel(column, fontsize = 18)
                  plt.ylabel('capacity_mw', fontsize = 18)
              plotnumber += 1
          plt.tight_layout()
```

```
------------------------
Scatter Plot :-
------------------------
```

**Features are related to label**

# Data Scaling

```
In [52]: scaler = StandardScaler()
         x_scaled = scaler.fit_transform(x)
         x_scaled
```

```
Out[52]: array([[ 1.15504214, -0.90475257,  1.1913948 , ...,  0.20019784,
                   0.1919957 ,  0.1785893 ],
                 [ 0.5985875 , -0.57860276, -0.98044461, ...,  0.20019784,
                   0.1919957 ,  0.1785893 ],
                 [ 0.13251441, -1.82661503,  1.62576269, ...,  0.20019784,
                   0.1919957 ,  0.1785893 ],
                 ...,
                 [-0.9466585 , -0.34689171,  1.62576269, ...,  0.20019784,
                   0.1919957 ,  0.1785893 ],
                 [ 0.53080541, -0.78390338, -0.98044461, ...,  0.20019784,
                   0.1919957 ,  0.1785893 ],
                 [-1.81634681,  0.10496987,  1.62576269, ...,  0.20019784,
                   0.1919957 ,  0.1785893 ]])
```

**Data has been scaled**

**Split data into train and test. Model will be bulit on training data and tested on test data**

```
In [53]: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.25, rand
         print('Data has been splited.')
```

```
Data has been splited.
```

# Model Bulding

### Decision Tree model instantiaing, training and evaluating

```
In [54]: DT = DecisionTreeRegressor()
         DT.fit(x_train, y_train)
         y_pred = DT.predict(x_test)
```

```
In [57]: print('=========================================================')
         print('R2 Score ----->', r2_score(y_test, y_pred))
         print('=========================================================')
         print('RMSE of Model ------->', np.sqrt(mean_squared_error(y_test, y_pred)))
         print('=========================================================')
         print('MSE of Model ------->', mean_squared_error(y_test, y_pred))
         print('=========================================================')
         print('Score of test data ---->', DT.score(x_test, y_test))
         print('=========================================================')
```

```
=========================================================
R2 Score -----> 0.6111258730830611
=========================================================
RMSE of Model -------> 274.22435732040987
=========================================================
MSE of Model -------> 75198.99814779183
=========================================================
Score of test data ----> 0.6111258730830611
=========================================================
```

**Conclusion : Decision Tree model has 61% score**

## XGBoost model instantiaing, training and evaluating

```
In [58]: xgb = xgb.XGBRegressor(eval_metric = 'mlogloss')
         xgb.fit(x_train, y_train)
         y_pred = xgb.predict(x_test)
```

```
In [59]: print('=========================================================')
         print('R2 Score ----->', r2_score(y_test, y_pred))
         print('=========================================================')
         print('RMSE of Model ------->', np.sqrt(mean_squared_error(y_test, y_pred)))
         print('=========================================================')
         print('MSE of Model ------->', mean_squared_error(y_test, y_pred))
         print('=========================================================')
         print('Score of test data ---->', xgb.score(x_test, y_test))
         print('=========================================================')
```

```
=========================================================
R2 Score -----> 0.752973794775826
=========================================================
RMSE of Model -------> 218.56117471317214
=========================================================
MSE of Model -------> 47768.98709200176
=========================================================
Score of test data ----> 0.752973794775826
=========================================================
```

**Conclusion : XGBoost model has 75% score**

## Knn model instantiaing, training and evaluating

```
In [60]: Knn = KNeighborsRegressor()
         Knn.fit(x_train, y_train)
         y_pred = Knn.predict(x_test)
```

```
In [61]: print('========================================================')
         print('R2 Score ----->', r2_score(y_test, y_pred))
         print('========================================================')
         print('RMSE of Model ------->', np.sqrt(mean_squared_error(y_test, y_pred)))
         print('========================================================')
         print('MSE of Model ------->', mean_squared_error(y_test, y_pred))
         print('========================================================')
         print('Score of test data ---->', Knn.score(x_test, y_test))
         print('========================================================')
```

```
========================================================
R2 Score -----> 0.7713684825660243
========================================================
RMSE of Model -------> 210.2662417010793
========================================================
MSE of Model -------> 44211.8923990967
========================================================
Score of test data ----> 0.7713684825660243
========================================================
```

**Conclusion : KNN model has 77% score**

## Random Forest model instantiaing, training and evaluating

```
In [62]: Rn = RandomForestRegressor()
         Rn.fit(x_train, y_train)
         y_pred = Rn.predict(x_test)
```

```
In [63]:  print('=========================================================')
          print('R2 Score ----->', r2_score(y_test, y_pred))
          print('=========================================================')
          print('RMSE of Model ------->', np.sqrt(mean_squared_error(y_test, y_pred)))
          print('=========================================================')
          print('MSE of Model ------->', mean_squared_error(y_test, y_pred))
          print('=========================================================')
          print('Score of test data ---->', Rn.score(x_test, y_test))
          print('=========================================================')
```

```
=========================================================
R2 Score -----> 0.8146813211243855
=========================================================
RMSE of Model -------> 189.30457093961715
=========================================================
MSE of Model -------> 35836.22057863254
=========================================================
Score of test data ----> 0.8146813211243855
=========================================================
```

**Conclusion : Random Forest model has 81% score**

## Looking RSME score we found Random Forest has best model so we do Hyperparameter Tuning on it

```
In [64]:  param_grid = {'n_estimators': [100,200,300,400,500],
                        'max_features': ['auto', 'squrt'],
                        'max_depth': [5, 10, 15, 20, 25, 30],
                        'min_samples_split': [2, 5, 10, 15, 100],
                        'min_samples_leaf': [1, 2, 5, 10]}
```

```
In [65]:  grid_search = GridSearchCV(estimator = Rn, param_grid = param_grid, cv = 5,n_jobs
```

```
In [66]:  grid_search.fit(x_train, y_train)
```

```
Out[66]:  GridSearchCV(cv=5, estimator=RandomForestRegressor(), n_jobs=-1,
                       param_grid={'max_depth': [5, 10, 15, 20, 25, 30],
                                   'max_features': ['auto', 'squrt'],
                                   'min_samples_leaf': [1, 2, 5, 10],
                                   'min_samples_split': [2, 5, 10, 15, 100],
                                   'n_estimators': [100, 200, 300, 400, 500]})
```

```
In [67]:  best_parameters = grid_search.best_params_
          print(best_parameters)
```

```
{'max_depth': 30, 'max_features': 'auto', 'min_samples_leaf': 5, 'min_samples_s
plit': 5, 'n_estimators': 200}
```

```
In [68]: hRn = RandomForestRegressor(max_depth = 30, max_features = 'auto', min_samples_le
         hRn.fit(x_train, y_train)
         hRn.score(x_test, y_test)
```

Out[68]: 0.8034611402413743

```
In [69]: y_pred = hRn.predict(x_test)
```

```
In [70]: print('=======================================================')
         print('R2 Score ----->', r2_score(y_test, y_pred))
         print('=======================================================')
         print('RMSE of Model ------->', np.sqrt(mean_squared_error(y_test, y_pred)))
         print('=======================================================')
         print('MSE of Model ------->', mean_squared_error(y_test, y_pred))
         print('=======================================================')
         print('Score of test data ---->', hRn.score(x_test, y_test))
         print('=======================================================')
```

```
=======================================================
R2 Score -----> 0.8034611402413743
=======================================================
RMSE of Model -------> 194.95111295013075
=======================================================
MSE of Model -------> 38005.93644049463
=======================================================
Score of test data ----> 0.8034611402413743
=======================================================
```

**After Hyperparameter Tuning model accuracy score 80%.**

## Saving The Model

```
In [71]: # saving the model to the Local file system
         filename = 'Global Power Plant Project (Capacity MW).pickle'
         pickle.dump(hRn, open(filename, 'wb'))
```

# Final Conclusion : Random Forest is our best model.

```
In [ ]:
```