

In [3]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import plotly.express as px
6 import plotly.graph_objects as go
7 import folium
8 from scipy.stats import pearsonr
9 from folium import plugins
10 from folium.plugins import HeatMap
11 from folium.plugins import HeatMapWithTime
12 import calendar
13 from matplotlib.ticker import FuncFormatter
14 import datetime
15 import joblib
16 import xgboost as xgb
17 import sklearn
18 from sklearn.model_selection import train_test_split
19 from sklearn.metrics import mean_squared_error
20 from sklearn.metrics import mean_absolute_error
21 from sklearn.model_selection import RandomizedSearchCV
22 from sklearn.model_selection import train_test_split
23 import statsmodels
24 from statsmodels.tsa.stattools import adfuller
25 from statsmodels.tsa.seasonal import seasonal_decompose
26 from statsmodels.tsa.arima_model import ARIMA
27 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
```

In [4]:

```
1 data = pd.read_excel('consumption.xlsx')
2 data.head()
```

Out[4]:

	Month	Year	Day of the Month	Rain fall,Inches	TotalTreated water leaving plant,MG	Hours Plant in operation	Total Raw water to plant,MG	Compliance with permitted capacity?	Backwas,Thousand Gallons	Peak demand into distribution,MDG	...	Measuremnt Recorded ≥ Measuremnt Required?	No Specified Treatn
0	January	2017	1	0	15.38	24	27.14	yes	3043	15.29	...	Yes	
1	January	2017	2	0	15.56	24	27.13	yes	3076	15.47	...	Yes	
2	January	2017	3	0	15.74	24	27.14	yes	3110	15.6	...	Yes	
3	January	2017	4	0	16.02	24	26.83	yes	3095	15.62	...	Yes	
4	January	2017	5	0	15.25	24	26.63	yes	3038	15.5	...	Yes	

5 rows × 22 columns

In [5]:

```
1 data1 = pd.read_csv('Cleaned.csv')
2 data1.head()
```

Out[5]:

	Date	Adjusted Daily Consumption, Kwh
0	01/01/2017	153910
1	02/01/2017	153910
2	03/01/2017	153910
3	04/01/2017	153910
4	05/01/2017	153910

In [6]:

```

1 column_names = data.columns.tolist()
2
3 # Print the List of column names
4 print("List of columns in the DataFrame:")
5 print(column_names)

```

List of columns in the DataFrame:

```

['Month', 'Year', 'Day of the Month', 'Rain fall,Inches', 'TotalTreated water leaving plant,MG', 'Hours Plant in operation', 'Total Raw water to plant,MG', 'Compliance with permitted capacity?', 'Backwas,Thousand Gallons', 'Peak demand into distribution,MDG', 'Numeasurements Recorded', 'No Measuremnts Required', 'Measuremnt Recorded ≥ Measuremnt Required?', 'N Measuremnt ≤ SpecifiedTreatmentTeachique Limit', 'N Measuremnt > Never Exceed Limit', 'Average Daily Turbidity,NTU', 'Maximum DailyTurbidity,NTU', 'Log Inactivation,Giardia', 'Giardia Compliance?', 'Log Inactivation,Viruses', 'Virus Compliance?', 'Emergency or Abnormal Operating Conditions,Repairor Maintenance Work that InvolvesTaking Water System Components Out of Operation']

```

In [7]:

```

1 # List of columns to check for non-numeric values
2 columns_to_check = ['Rain fall,Inches', 'TotalTreated water leaving plant,MG', 'Hours Plant in operation', 'Total Raw water to plant,MG', 'Compliance with permitted capacity?', 'Backwas,Thousand Gallons', 'Peak demand into distribution,MDG', 'Numeasurements Recorded', 'No Measuremnts Required', 'Average Daily Turbidity,NTU', 'Maximum DailyTurbidity,NTU', 'Log Inactivation,Giardia', 'Log Inactivation,Viruses']
3
4 # Convert to numeric, coercing non-numeric to NaN
5 for col in columns_to_check:
6     data[col] = pd.to_numeric(data[col], errors='coerce')
7
8 # Drop rows where any of the specified columns are NaN
9 data.dropna(subset=columns_to_check, inplace=True)
10
11 # Reset the index
12 data.reset_index(drop=True, inplace=True)

```

In [8]:

```

1 column_data_types = data.dtypes
2
3 # Print the data types
4 print("Data types of each column in the DataFrame:")
5 print(column_data_types)

```

Data types of each column in the DataFrame:

```

Month
object
Year
int64
Day of the Month
object
Rain fall,Inches
float64
TotalTreated water leaving plant,MG
float64
Hours Plant in operation
float64
Total Raw water to plant,MG
float64
Compliance with permitted capacity?
object
Backwas,Thousand Gallons
float64
Peak demand into distribution,MDG
float64
Numeasurements Recorded
float64
No Measuremnts Required
float64
Measuremnt Recorded ≥ Measuremnt Required?
object
N Measuremnt ≤ SpecifiedTreatmentTeachique Limit
object
N Measuremnt > Never Exceed Limit
object
Average Daily Turbidity,NTU
float64
Maximum DailyTurbidity,NTU
float64
Log Inactivation,Giardia
float64
Giardia Compliance?
object
Log Inactivation,Viruses
float64
Virus Compliance?
object
Emergency or Abnormal Operating Conditions,Repairor Maintenance Work that InvolvesTaking Water System Components Out of Operation
object
dtype: object

```

In [9]:

```

1 # Correct the year naming issue (changing "20222" to "2022")
2 data['Year'] = data['Year'].replace(20222, 2022)
3
4 # Convert Month, Day, and Year columns to a datetime format
5 data['Date'] = pd.to_datetime(data['Year'].astype(str) + '-' + data['Month'].astype(str) + '-' + data['Day of the Month'].astype(str))
6

```

In [10]:

```

1 # Extract the required columns
2 data = data[['Date', 'Year', 'Month', 'Day of the Month', 'Rain fall,Inches', 'TotalTreated water leaving plant,MG', 'Hours Plant i
3             'Backwas,Thousand Gallons', 'Peak demand into distribution,MDG', 'Nomeasurements Recorded',
4             'No Measuremnts Required', 'Average Daily Turbidity,NTU', 'Maximum DailyTurbidity,NTU', 'Log Inactivation,Giardi
5             'Log Inactivation,Viruses']]
6
7 # Display the extracted data
8 data.head()

```

Out[10]:

	Date	Year	Month	Day of the Month	Rain fall,Inches	TotalTreated water leaving plant,MG	Hours Plant in operation	Total Raw water to plant,MG	Backwas,Thousand Gallons	Peak demand into distribution,MDG	Nomeasurements Recorded	No Measuremnts Required	Turbi
0	2017-01-01	2017	January	1	0.0	15.38	24.0	27.14	3043.0	15.29	6.0	6.0	
1	2017-01-02	2017	January	2	0.0	15.56	24.0	27.13	3076.0	15.47	6.0	6.0	
2	2017-01-03	2017	January	3	0.0	15.74	24.0	27.14	3110.0	15.60	6.0	6.0	
3	2017-01-04	2017	January	4	0.0	16.02	24.0	26.83	3095.0	15.62	6.0	6.0	
4	2017-01-05	2017	January	5	0.0	15.25	24.0	26.63	3038.0	15.50	6.0	6.0	

In [11]:

```

1 # Convert the 'Date' column in data1 to datetime format
2 data1['Date'] = pd.to_datetime(data1['Date'])
3
4 # Now you can perform the merge
5 data = pd.merge(data, data1, on='Date', how='inner')
6

```

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '13/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '14/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '15/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '16/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '17/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '18/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

C:\Users\User\AppData\Local\Temp\ipykernel_18284\3107944606.py:2: UserWarning: Parsing '19/01/2017' in DD/MM/YYYY format. Provide format or specify infer_datetime_format=True for consistent parsing.
data1['Date'] = pd.to_datetime(data1['Date'])

In [12]:

```

1 column_data_types = data.dtypes
2
3 # Print the data types
4 print("Data types of each column in the DataFrame:")
5 print(column_data_types)

```

Data types of each column in the DataFrame:

```

Date                                datetime64[ns]
Year                                int64
Month                                object
Day of the Month                     object
Rain fall,Inches                     float64
TotalTreated water leaving plant,MG  float64
Hours Plant in operation              float64
Total Raw water to plant,MG          float64
Backwas,Thousand Gallons             float64
Peak demand into distribution,MDG    float64
Nomeasurements Recorded              float64
No Measurementns Required             float64
Average Daily Turbidity,NTU          float64
Maximum DailyTurbidity,NTU           float64
Log Inactivation,Giardia              float64
Log Inactivation,Viruses              float64
Adjusted Daily Consumption, Kwh       int64
dtype: object

```

In [14]:

```

1 data = pd.get_dummies(data, columns=['Month'], drop_first=True)
2
3 # Convert 'Day of the Month' to integer if it's an object
4 if data['Day of the Month'].dtype == 'object':
5     data['Day of the Month'] = data['Day of the Month'].astype(int)
6

```

In [15]:

```

1 # Check the data types of each column
2 print('Data types of each column in the DataFrame:')
3 print(data.dtypes)
4
5 # Check for missing values
6 print('\nMissing values in each column:')
7 print(data.isnull().sum())
8
9 # Display basic statistics
10 print('\nBasic statistics:')
11 print(data.describe())
12

```

count	877.000000	877.000000	877.000000	877.000000	877.000000
mean	0.119726	0.045610	0.099202	0.196123	0.011403
std	0.324826	0.208757	0.299103	0.397289	0.106233
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.000000	0.000000	0.000000	0.000000	0.000000
max	1.000000	1.000000	1.000000	1.000000	1.000000

	Month_June	Month_March	Month_May	Month_November
count	877.000000	877.000000	877.000000	877.000000
mean	0.051311	0.174458	0.119726	0.027366
std	0.220758	0.379720	0.324826	0.163241
min	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000
50%	0.000000	0.000000	0.000000	0.000000
75%	0.000000	0.000000	0.000000	0.000000
max	1.000000	1.000000	1.000000	1.000000

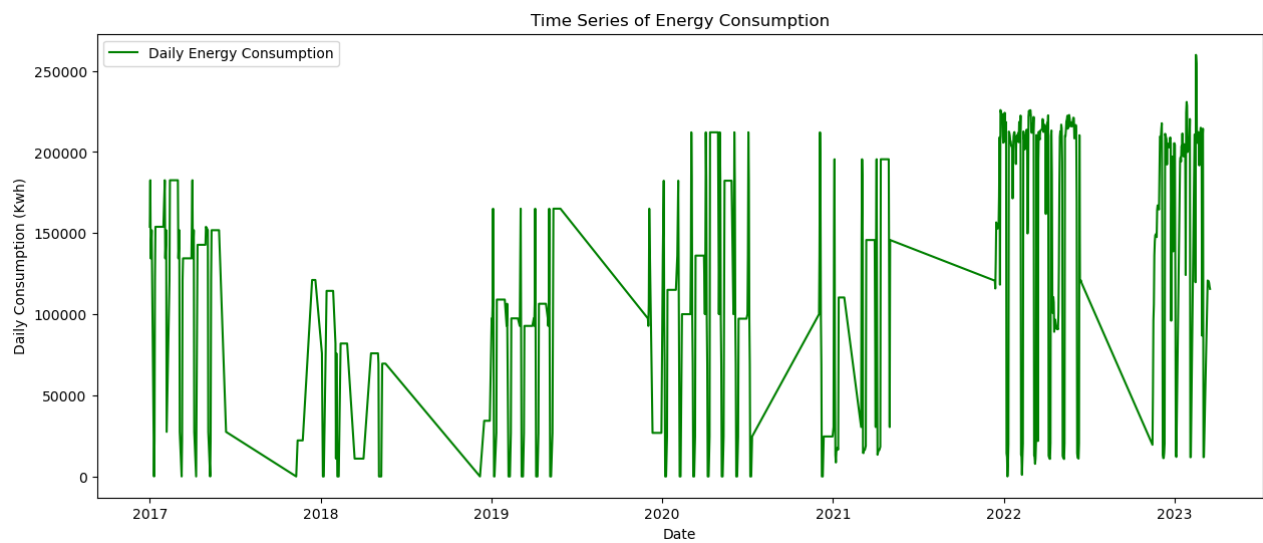
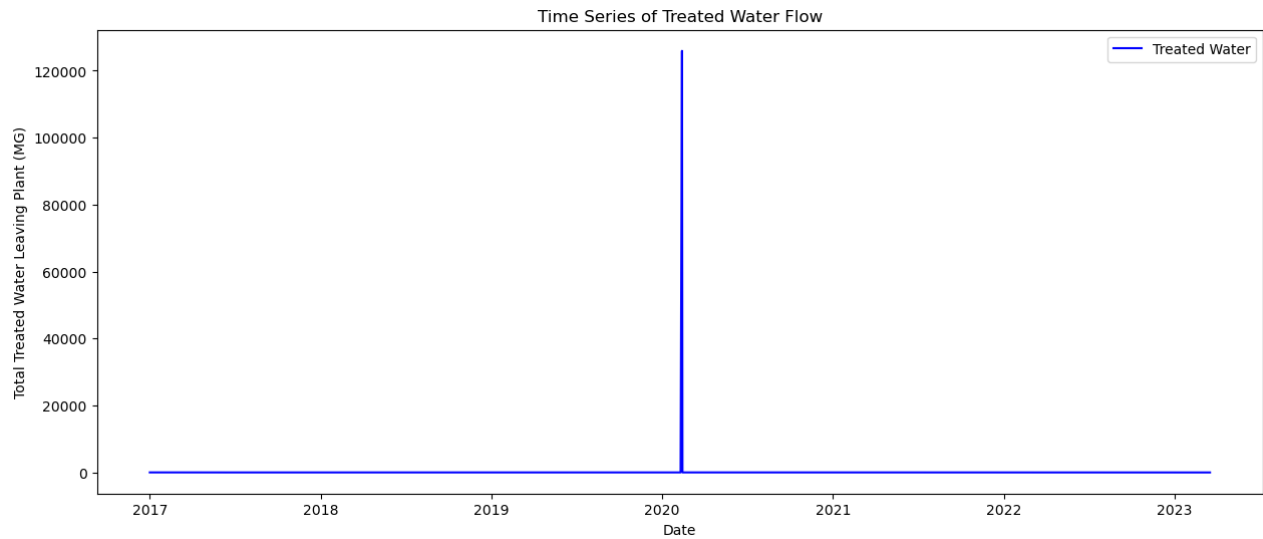
[8 rows x 24 columns]

In [16]:

```

1 # Plot for Total Treated Water Leaving Plant
2 plt.figure(figsize=(15, 6))
3 plt.plot(data['Date'], data['TotalTreated water leaving plant,MG'], label='Treated Water', color='blue')
4 plt.xlabel('Date')
5 plt.ylabel('Total Treated Water Leaving Plant (MG)')
6 plt.title('Time Series of Treated Water Flow')
7 plt.legend()
8 plt.show()
9
10 # Plot for Total Raw Water to Plant
11 plt.figure(figsize=(15, 6))
12 plt.plot(data['Date'], data['Adjusted Daily Consumption, Kwh'], label='Daily Energy Consumption', color='green')
13 plt.xlabel('Date')
14 plt.ylabel('Daily Consumption (Kwh)')
15 plt.title('Time Series of Energy Consumption')
16 plt.legend()
17 plt.show()
18

```



In [17]:

```

1 # Removing outliers beyond 99th percentile
2 for col in ['TotalTreated water leaving plant,MG', 'Total Raw water to plant,MG', 'Peak demand into distribution,MDG']:
3     cap = data[col].quantile(0.99)
4     data = data[data[col] <= cap]

```

In [18]:

1 data.head()

Out[18]:

	Date	Year	Day of the Month	Rain fall, Inches	Total Treated water leaving plant, MG	Hours Plant in operation	Total Raw water to plant, MG	Backwas, Thousand Gallons	Peak demand into distribution, MDG	Noneasurements Recorded	...	Adjusted Daily Consumption, Kwh	Month_De
0	2017-01-01	2017	1	0.0	15.38	24.0	27.14	3043.0	15.29	6.0	...	153910	
1	2017-01-02	2017	2	0.0	15.56	24.0	27.13	3076.0	15.47	6.0	...	182614	
2	2017-01-03	2017	3	0.0	15.74	24.0	27.14	3110.0	15.60	6.0	...	134468	
3	2017-01-04	2017	4	0.0	16.02	24.0	26.83	3095.0	15.62	6.0	...	142808	
4	2017-01-05	2017	5	0.0	15.25	24.0	26.63	3038.0	15.50	6.0	...	151815	

5 rows × 25 columns

In [19]:

```

1 # Checking for non-numeric values
2
3 # Function to check if a string contains only digits and periods
4 def check_string(value):
5     for char in str(value):
6         if char not in '1234567890.':
7             return False
8     return True
9
10 # Dictionary to store columns with non-integer or non-period characters
11 columns_with_non_integers = {}
12
13 # Loop through each column that is supposed to be numeric after preprocessing
14 for col in data: # Replace with your actual numeric columns
15     if col != "Date":
16         # Check each value in the column
17         for value in data[col]:
18             if not check_string(value):
19                 if col not in columns_with_non_integers:
20                     columns_with_non_integers[col] = []
21                     columns_with_non_integers[col].append(value)
22
23 # Display the columns and their non-integer values
24 for col, values in columns_with_non_integers.items():
25     print(f'Column {col} has non-integer values: {values[:10]}...') # Display only the first 10 non-integer values

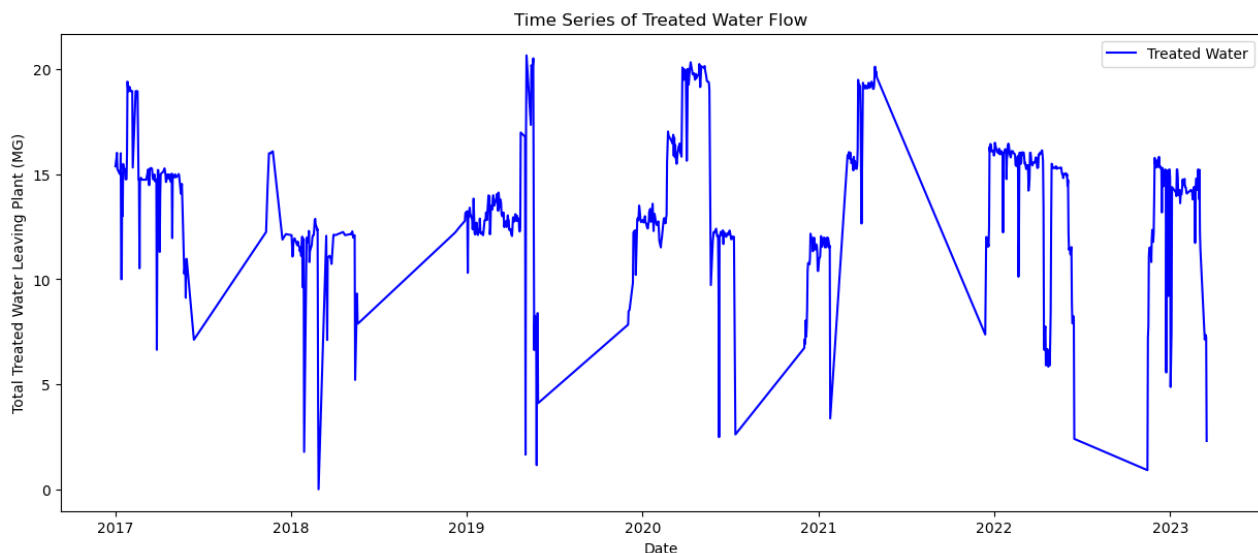
```

In [20]:

```

1 # Plot for Total Treated Water Leaving Plant
2 plt.figure(figsize=(15, 6))
3 plt.plot(data['Date'], data['TotalTreated water leaving plant, MG'], label='Treated Water', color='blue')
4 plt.xlabel('Date')
5 plt.ylabel('Total Treated Water Leaving Plant (MG)')
6 plt.title('Time Series of Treated Water Flow')
7 plt.legend()
8 plt.show()
9

```

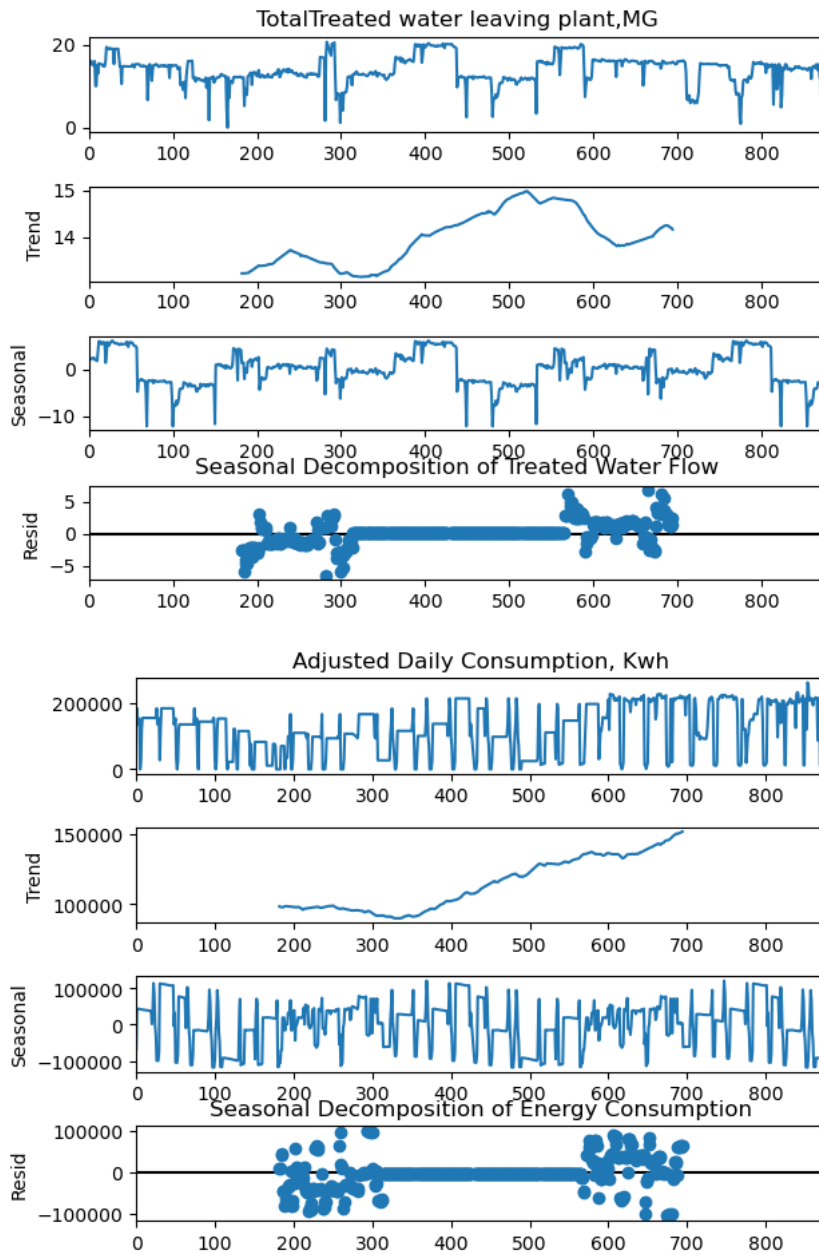


In [21]:

```

1 # Seasonal Decomposition for Treated Water
2 decomposition_treated = seasonal_decompose(data['TotalTreated water leaving plant,MG'], period=365)
3 decomposition_treated.plot()
4 plt.title('Seasonal Decomposition of Treated Water Flow')
5 plt.show()
6
7 # Seasonal Decomposition for Raw Water
8 decomposition_raw = seasonal_decompose(data['Adjusted Daily Consumption, Kwh'], period=365)
9 decomposition_raw.plot()
10 plt.title('Seasonal Decomposition of Energy Consumption')
11 plt.show()
12

```



In [22]:

```
1 from statsmodels.tsa.seasonal import seasonal_decompose
2
3 # Perform seasonal decomposition for Treated Water
4
5 treated_seasonal = decomposition_treated.seasonal
6 treated_trend = decomposition_treated.trend
7 treated_residual = decomposition_treated.resid
8
9 # Perform seasonal decomposition for Raw Water
10 raw_seasonal = decomposition_raw.seasonal
11 raw_trend = decomposition_raw.trend
12 raw_residual = decomposition_raw.resid
13
14 # Create DataFrames to hold the decomposition components for each target variable
15 decomposition_treated_df = pd.DataFrame({
16     'Date': data['Date'],
17     'Seasonal': treated_seasonal,
18     'Trend': treated_trend,
19     'Residual': treated_residual
20 })
21
22 decomposition_raw_df = pd.DataFrame({
23     'Date': data['Date'],
24     'Seasonal': raw_seasonal,
25     'Trend': raw_trend,
26     'Residual': raw_residual
27 })
28
29 # Display the first few rows of each DataFrame
30 print("Decomposition Components for Treated Water:")
31 print(decomposition_treated_df.head())
32
33 print("\nDecomposition Components for Energy Consumption:")
34 print(decomposition_raw_df.head())
35
36
37
38 # Create DataFrames to hold the decomposition components for each target variable
39 decomposition_treated_df = pd.DataFrame({
40     'Date': data['Date'],
41     'Seasonal': treated_seasonal,
42     'Trend': treated_trend,
43     'Residual': treated_residual
44 })
45
46 decomposition_raw_df = pd.DataFrame({
47     'Date': data['Date'],
48     'Seasonal': raw_seasonal,
49     'Trend': raw_trend,
50     'Residual': raw_residual
51 })
52
53 # Display the first few non-NaN rows of each DataFrame
54 print("Decomposition Components for Treated Water:")
55 print(decomposition_treated_df.dropna().head())
56
57 print("\nDecomposition Components for Daily Consumption, Kwh:")
58 print(decomposition_raw_df.dropna().head())
59
60
```


Decomposition Components for Treated Water:

	Date	Seasonal	Trend	Residual
0	2017-01-01	1.664258	NaN	NaN
1	2017-01-02	1.886807	NaN	NaN
2	2017-01-03	2.157730	NaN	NaN
3	2017-01-04	2.225113	NaN	NaN
4	2017-01-05	2.376322	NaN	NaN

Decomposition Components for Energy Consumption:

	Date	Seasonal	Trend	Residual
0	2017-01-01	-69051.966857	NaN	NaN
1	2017-01-02	42405.934513	NaN	NaN
2	2017-01-03	42551.756431	NaN	NaN
3	2017-01-04	42692.142732	NaN	NaN
4	2017-01-05	42346.934513	NaN	NaN

Decomposition Components for Treated Water:

	Date	Seasonal	Trend	Residual
182	2018-05-10	1.345421	13.227793	-2.593213
183	2018-05-11	1.428341	13.228119	-2.606460
184	2018-05-13	1.392558	13.228615	-2.521173
185	2018-05-14	-2.058347	13.229413	-5.961067
186	2018-05-18	0.132310	13.229733	-4.042043

Decomposition Components for Daily Consumption, Kwh:

	Date	Seasonal	Trend	Residual
182	2018-05-10	-105600.116172	98094.030137	7506.086035
183	2018-05-11	-105532.798364	97739.865753	7792.932610
184	2018-05-13	-70695.391514	97612.545205	42627.846309
185	2018-05-14	-69819.966857	97617.131507	41747.835350
186	2018-05-18	18693.646842	97598.868493	-46747.515335

In [23]:

```
1 from statsmodels.tsa.stattools import adfuller
2
3 # Perform Augmented Dickey-Fuller test
4 def adf_test(series):
5     result = adfuller(series, autolag='AIC')
6     print(f'ADF Statistic: {result[0]}')
7     print(f'p-value: {result[1]}')
8     print(f'Critical Values: {result[4]}')
9     if result[1] <= 0.05:
10         print('Strong evidence against the null hypothesis, reject the null hypothesis. Data has no unit root and is stationary')
11     else:
12         print('Weak evidence against null hypothesis, time series has a unit root, indicating it is non-stationary ')
13
14 # Apply the test on your time series data
15 adf_test(data['TotalTreated water leaving plant,MG'])
16 adf_test(data['Adjusted Daily Consumption, Kwh'])
17
```

ADF Statistic: -4.533456769576553

p-value: 0.00017100684509673158

Critical Values: {'1%': -3.4381216826257956, '5%': -2.8649705364894635, '10%': -2.568596692178972}

Strong evidence against the null hypothesis, reject the null hypothesis. Data has no unit root and is stationary

ADF Statistic: -6.58512382505417

p-value: 7.3474096452009324e-09

Critical Values: {'1%': -3.4381032536542913, '5%': -2.8649624121419746, '10%': -2.5685923644574107}

Strong evidence against the null hypothesis, reject the null hypothesis. Data has no unit root and is stationary

In [24]:

```
1 from statsmodels.tsa.statespace.sarimax import SARIMAX
2 from sklearn.metrics import mean_squared_error, mean_absolute_error
3 from math import sqrt
4
5 # Train-Test Split
6 data = data.sort_values('Date')
7 train_size = int(len(data) * 0.8)
8 train, test = data[0:train_size], data[train_size:len(data)]
```

In [25]:

```

1 # SARIMA Model Training for Treated Water
2 sarima_model_treated = SARIMAX(train['TotalTreated water leaving plant,MG'],
3                                order=(1, 1, 1),
4                                seasonal_order=(1, 1, 1, 12))
5 sarima_fit_treated = sarima_model_treated.fit(dispatch=False)
6
7 # SARIMA Model Training for Raw Water
8 sarima_model_raw = SARIMAX(train['Adjusted Daily Consumption, Kwh'],
9                             order=(1, 1, 1),
10                             seasonal_order=(1, 1, 1, 12))
11 sarima_fit_raw = sarima_model_raw.fit(dispatch=False)
12

```

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

In [26]:

```

1 # Model Evaluation
2 # SARIMA
3 sarima_pred_treated = sarima_fit_treated.predict(len(train), len(data)-1)
4 sarima_pred_raw = sarima_fit_raw.predict(len(train), len(data)-1)

```

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:834: ValueWarning: No supported index is available. Prediction results will be given with an integer index beginning at 'start'.

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:834: ValueWarning: No supported index is available. Prediction results will be given with an integer index beginning at 'start'.

In [27]:

```

1 # Calculate RMSE and MAE
2 print("SARIMA RMSE for Treated Water:", sqrt(mean_squared_error(test['TotalTreated water leaving plant,MG'], sarima_pred_treated)))
3 print("SARIMA RMSE for Energy Consumption:", sqrt(mean_squared_error(test['Adjusted Daily Consumption, Kwh'], sarima_pred_raw)))

```

SARIMA RMSE for Treated Water: 4.7015497521295115

SARIMA RMSE for Energy Consumption: 72614.11760287227

In [28]:

```

1 print("SARIMA MAE for Treated Water:", mean_absolute_error(test['TotalTreated water leaving plant,MG'], sarima_pred_treated))
2 print("SARIMA MAE for Energy Consumption:", mean_absolute_error(test['Adjusted Daily Consumption, Kwh'], sarima_pred_raw))
3
4 #print("Prophet MAE for Treated Water:", mean_absolute_error(test['TotalTreated water leaving plant,MG'], prophet_pred_treated))
5 #print("Prophet MAE for Raw Water:", mean_absolute_error(test['Total Raw water to plant,MG'], prophet_pred_raw))

```

SARIMA MAE for Treated Water: 3.2650915054531637

SARIMA MAE for Energy Consumption: 53888.83678841617

Random Forest as the machine learning model and SARIMA as the time series model. It then averages the predictions from both models to create an ensemble. You can adjust the ensemble weights based on the performance of individual models.

In [29]:

```

1 # Importing required Libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
5 from sklearn.metrics import mean_squared_error, mean_absolute_error
6 from sklearn.model_selection import train_test_split
7 from statsmodels.tsa.statespace.sarimax import SARIMAX
8 from math import sqrt

```

In [30]:

```
1 # Assuming 'data' is your DataFrame containing historical data
2 # Log transformation of the target variables
3 data['Log_TotalTreated'] = np.log1p(data['TotalTreated water leaving plant,MG'])
4 data['Log_EnergyConsumption'] = np.log1p(data['Adjusted Daily Consumption, Kwh'])
5
6 # Feature Engineering: Create lag features for target variables
7 for col in ['Log_EnergyConsumption']:
8     for lag in range(1, 4): # Create 3 Lag features
9         data[f'{col}_lag{lag}'] = data[col].shift(lag)
10
11 # Drop NaN values created due to lag features
12 data.dropna(inplace=True)
```

In [122]:

```
1 # Split sizes
2 train_size = int(len(data) * 0.6) # 70% for training
3 validation_size = int(len(data) * 0.2) # 20% for validation
4
5 # Train-Validation-Test Split
6 train, temp = data[:train_size], data[train_size:]
7 validation, test = temp[:validation_size], temp[validation_size:]
8
9 # Prepare data for machine learning models for actual values
10
11 # Training data
12 X_train_actual = train.drop(['Date', 'TotalTreated water leaving plant,MG', 'Adjusted Daily Consumption, Kwh'], axis=1)
13 y_train_treated_actual = train['TotalTreated water leaving plant,MG']
14 y_train_energy_actual = train['Adjusted Daily Consumption, Kwh']
15
16 # Validation data
17 X_validation_actual = validation.drop(['Date', 'TotalTreated water leaving plant,MG', 'Adjusted Daily Consumption, Kwh'], axis=1)
18 y_validation_treated_actual = validation['TotalTreated water leaving plant,MG']
19 y_validation_energy_actual = validation['Adjusted Daily Consumption, Kwh']
20
21 # Test data
22 X_test_actual = test.drop(['Date', 'TotalTreated water leaving plant,MG', 'Adjusted Daily Consumption, Kwh'], axis=1)
23 y_test_treated_actual = test['TotalTreated water leaving plant,MG']
24 y_test_energy_actual = test['Adjusted Daily Consumption, Kwh']
25
26 print("Training set size:", len(X_train_actual))
27 print("Validation set size:", len(X_validation_actual))
28 print("Test set size:", len(X_test_actual))
29
```

Training set size: 510
Validation set size: 170
Test set size: 170

In [68]:

```
1 X_train_actual
```

Out[68]:

	Year	Day of the Month	Rain fall,Inches	Hours Plant in operation	Total Raw water to plant,MG	Backwas,Thousand Gallons	Peak demand into distribution,MDG	Nomeasurements Recorded	No Measuremnts Required	Average Daily Turbidity,NTU	...	Month_July
3	2017	4	0.0	24.0	26.830000	3095.0000	15.6200	6.0	6.0	0.0200	...	0
4	2017	5	0.0	24.0	26.630000	3038.0000	15.5000	6.0	6.0	0.0200	...	0
5	2017	10	0.0	24.0	26.900000	3045.0000	15.6400	6.0	6.0	0.0200	...	0
6	2017	11	0.0	24.0	26.530000	2982.0000	15.5700	6.0	6.0	0.0000	...	0
7	2017	12	0.0	24.0	26.320000	2951.0000	15.3300	6.0	6.0	0.0200	...	0
...
616	2022	8	0.0	24.0	29.462891	3245.9990	16.2395	6.0	6.0	0.0722	...	0
617	2022	9	0.0	24.0	18.392578	3311.4417	16.2369	6.0	6.0	0.0722	...	0
618	2022	10	0.0	24.0	20.433594	3280.6940	16.2481	6.0	6.0	0.0719	...	0
619	2022	11	0.0	24.0	22.681641	3158.2337	19.4189	6.0	6.0	0.0723	...	0
620	2022	12	0.0	24.0	17.880859	3193.6305	16.2808	6.0	6.0	0.0756	...	0

595 rows × 27 columns

In [150]:

```

1
2 # Machine Learning Model Training for Treated Water
3 rf_model_treated = RandomForestRegressor()
4 rf_model_treated.fit(X_train_actual, y_train_treated_actual)
5
6 gb_model_treated = GradientBoostingRegressor()
7 gb_model_treated.fit(X_train_actual, y_train_treated_actual)
8
9 # Machine Learning Model Training for Energy Consumption
10 rf_model_energy = RandomForestRegressor()
11 rf_model_energy.fit(X_train_actual, y_train_energy_actual)
12
13 gb_model_energy = GradientBoostingRegressor()
14 gb_model_energy.fit(X_train_actual, y_train_energy_actual)
15
16 # SARIMA Model Training for Treated Water
17 try:
18     sarima_model_treated = SARIMAX(y_train_treated_actual, order=(1, 1, 1), seasonal_order=(1, 1, 1, 12))
19     sarima_fit_treated = sarima_model_treated.fit(dispatch=False)
20 except Exception as e:
21     print("Error fitting SARIMA model for Treated Water:", str(e))
22
23 # SARIMA Model Training for Energy Consumption
24 try:
25     sarima_model_energy = SARIMAX(y_train_energy_actual, order=(1, 1, 1), seasonal_order=(1, 1, 1, 12))
26     sarima_fit_energy = sarima_model_energy.fit(dispatch=False)
27 except Exception as e:
28     print("Error fitting SARIMA model for Energy Consumption:", str(e))

```

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

self._init_dates(dates, freq)
C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

self._init_dates(dates, freq)
C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

self._init_dates(dates, freq)
C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:471: ValueWarning: An unsupported index was provided and will be ignored when e.g. forecasting.

self._init_dates(dates, freq)

In [149]:

```

1 # Making Predictions for Actual Values on Validation Set
2 rf_pred_treated_validation = rf_model_treated.predict(X_validation_actual)
3 gb_pred_treated_validation = gb_model_treated.predict(X_validation_actual)
4 rf_pred_energy_validation = rf_model_energy.predict(X_validation_actual)
5 gb_pred_energy_validation = gb_model_energy.predict(X_validation_actual)
6
7 sarima_pred_treated_validation = sarima_fit_treated.predict(len(y_train_treated_actual), len(y_train_treated_actual) + len(y_validation_actual))
8 sarima_pred_energy_validation = sarima_fit_energy.predict(len(y_train_energy_actual), len(y_train_energy_actual) + len(y_validation_actual))
9
10 # Convert SARIMA predictions to numpy arrays for consistency
11 sarima_pred_treated_validation = np.array(sarima_pred_treated_validation)
12 sarima_pred_energy_validation = np.array(sarima_pred_energy_validation)
13
14 # Ensemble Predictions for Actual Values on Validation Set
15 treated_validation_outputs = [rf_pred_treated_validation, gb_pred_treated_validation, sarima_pred_treated_validation]
16 energy_validation_outputs = [rf_pred_energy_validation, gb_pred_energy_validation, sarima_pred_energy_validation]
17 ensemble_pred_treated_validation = np.mean(treated_validation_outputs, axis=0)
18 ensemble_pred_energy_validation = np.mean(energy_validation_outputs, axis=0)
19
20

```

C:\Users\User\anaconda3\lib\site-packages\statsmodels\tsa\base\tsa_model.py:834: ValueWarning: No supported index is available. Prediction results will be given with an integer index beginning at 'start'.

return get_prediction_index(

In [125]:

```

1 def error(y_true, y_pred, threshold=0.5):
2     """
3     Calculate the error rate based on a given threshold.
4     The function returns the proportion of predictions that are off by more than the threshold.
5     """
6     return np.mean(np.abs(y_true - y_pred) > threshold)

```

In [126]:

```

1 from sklearn.metrics import mean_squared_error, mean_absolute_error
2 from math import sqrt
3
4 def calculate_metrics(y_true, y_pred, model_name, target_name):
5     rmse = sqrt(mean_squared_error(y_true, y_pred))
6     mae = mean_absolute_error(y_true, y_pred)
7     error_rate = calculate_error_rate(y_true, y_pred) # Define this function!
8
9     print(f"=== {model_name} for {target_name} ===")
10    print(f"RMSE: {rmse:.2f}")
11    print(f"MAE: {mae:.2f}")
12    print(f"Error Rate: {error_rate:.2f}%\n")
13
14 # For example: calculate_error_rate can be a relative error, defined as:
15 def calculate_error_rate(y_true, y_pred):
16     return np.mean(np.abs((y_true - y_pred) / y_true)) * 100 # Mean Absolute Percentage Error (MAPE)
17
18 # Predictions and Model Names for Treated Water on Validation Set
19 treated_preds = [rf_pred_treated_validation, gb_pred_treated_validation, sarima_pred_treated_validation, ensemble_pred_treated_validation]
20 treated_model_names = ['RF', 'GB', 'SARIMA', 'Ensemble']
21
22 for model_name, pred in zip(treated_model_names, treated_preds):
23     calculate_metrics(y_validation_treated_actual, pred, model_name, 'Treated Water Validation')
24
25 # Predictions and Model Names for Energy Consumption on Validation Set
26 energy_preds = [rf_pred_energy_validation, gb_pred_energy_validation, sarima_pred_energy_validation, ensemble_pred_energy_validation]
27 energy_model_names = ['RF', 'GB', 'SARIMA', 'Ensemble']
28
29 for model_name, pred in zip(energy_model_names, energy_preds):
30     calculate_metrics(y_validation_energy_actual, pred, model_name, 'Energy Consumption Validation')

```

```

=== RF for Treated Water Validation ===
RMSE: 0.11
MAE: 0.05
Error Rate: 0.41%

```

```

=== GB for Treated Water Validation ===
RMSE: 0.07
MAE: 0.03
Error Rate: 0.32%

```

```

=== SARIMA for Treated Water Validation ===
RMSE: 5.30
MAE: 4.99
Error Rate: 31.39%

```

```

=== Ensemble for Treated Water Validation ===
RMSE: 1.77
MAE: 1.66
Error Rate: 10.50%

```

```

=== RF for Energy Consumption Validation ===
RMSE: 4876.59
MAE: 3422.66
Error Rate: 9.53%

```

```

=== GB for Energy Consumption Validation ===
RMSE: 4332.31
MAE: 2619.35
Error Rate: 13.41%

```

```

=== SARIMA for Energy Consumption Validation ===
RMSE: 54726.09
MAE: 30672.51
Error Rate: 1585.03%

```

```

=== Ensemble for Energy Consumption Validation ===
RMSE: 19062.70
MAE: 11461.67
Error Rate: 531.42%

```

In [127]:

```

1 from sklearn.model_selection import GridSearchCV
2
3 # Define the parameter grid
4 param_grid = {
5     'n_estimators': [50, 100, 200],
6     'learning_rate': [0.01, 0.1, 0.2],
7     'max_depth': [3, 4, 5],
8     'subsample': [0.8, 0.9, 1.0],
9     'min_samples_split': [2, 3, 4]
10 }
11
12 # Create the model to be tuned
13 gb_base = GradientBoostingRegressor()
14
15 # Create the grid search object
16 grid_search = GridSearchCV(estimator = gb_base, param_grid = param_grid,
17                             cv = 3, n_jobs = -1, verbose = 2)
18
19 # Fit the grid search to the data
20 grid_search.fit(X_train_actual, y_train_treated_actual)
21
22 # Get the best parameters
23 best_params = grid_search.best_params_
24 print("Best parameters:", best_params)
25
26 # Train and test the best model
27 best_gb_model_treated = GradientBoostingRegressor(**best_params)
28 best_gb_model_treated.fit(X_train_actual, y_train_treated_actual)

```

Fitting 3 folds for each of 243 candidates, totalling 729 fits

Best parameters: {'learning_rate': 0.1, 'max_depth': 3, 'min_samples_split': 4, 'n_estimators': 50, 'subsample': 0.9}

Out[127]:

GradientBoostingRegressor(min_samples_split=4, n_estimators=50, subsample=0.9)

In [128]:

```

1 from sklearn.model_selection import GridSearchCV
2
3 # Define the parameter grid
4 param_grid = {
5     'n_estimators': [50, 100, 200],
6     'learning_rate': [0.01, 0.1, 0.2],
7     'max_depth': [3, 4, 5],
8     'subsample': [0.8, 0.9, 1.0],
9     'min_samples_split': [2, 3, 4]
10 }
11
12 # Create the model to be tuned
13 gb_base = GradientBoostingRegressor()
14
15 # Create the grid search object
16 grid_search = GridSearchCV(estimator = gb_base, param_grid = param_grid,
17                             cv = 3, n_jobs = -1, verbose = 2)
18
19 # Fit the grid search to the data
20 grid_search.fit(X_train_actual, y_train_energy_actual)
21
22 # Get the best parameters
23 best_params = grid_search.best_params_
24 print("Best parameters:", best_params)
25
26 # Train and test the best model
27 best_gb_model_energy = GradientBoostingRegressor(**best_params)
28 best_gb_model_energy.fit(X_train_actual, y_train_energy_actual)

```

Fitting 3 folds for each of 243 candidates, totalling 729 fits

Best parameters: {'learning_rate': 0.1, 'max_depth': 3, 'min_samples_split': 3, 'n_estimators': 100, 'subsample': 0.9}

Out[128]:

GradientBoostingRegressor(min_samples_split=3, subsample=0.9)

In [129]:

```

1  # Making predictions on the training set
2  y_pred_train_treated = best_gb_model.predict(X_train_actual)
3
4  # Calculate RMSE, MAE, and Error Rate for the training set
5  rmse_train_treated = sqrt(mean_squared_error(y_train_treated_actual, y_pred_train_treated))
6  mae_train_treated = mean_absolute_error(y_train_treated_actual, y_pred_train_treated)
7  error_train_treated = error(y_train_treated_actual, y_pred_train_treated)
8
9  # Making predictions on the test set (as you already did)
10 y_pred_tuned_treated = best_gb_model.predict(X_test_actual)
11
12 # Calculate RMSE, MAE, and Error Rate for the test set
13 rmse_test_treated = sqrt(mean_squared_error(y_test_treated_actual, y_pred_tuned_treated))
14 mae_test_treated = mean_absolute_error(y_test_treated_actual, y_pred_tuned_treated)
15 error_test_treated = error(y_test_treated_actual, y_pred_tuned_treated)
16
17 # Displaying the results
18 print("=== Hyperparameter-tuned Model Evaluation for Treated Water ===")
19 print("\n--- Training Set ---")
20 print(f"RMSE: {rmse_train_treated}")
21 print(f"MAE: {mae_train_treated}")
22 print(f"Error Rate: {error_train_treated}")
23
24 print("\n--- Test Set ---")
25 print(f"RMSE: {rmse_test_treated}")
26 print(f"MAE: {mae_test_treated}")
27 print(f"Error Rate: {error_test_treated}")
28

```

=== Hyperparameter-tuned Model Evaluation for Treated Water ===

--- Training Set ---
 RMSE: 0.011638433832252987
 MAE: 0.008573316842963906
 Error Rate: 0.0

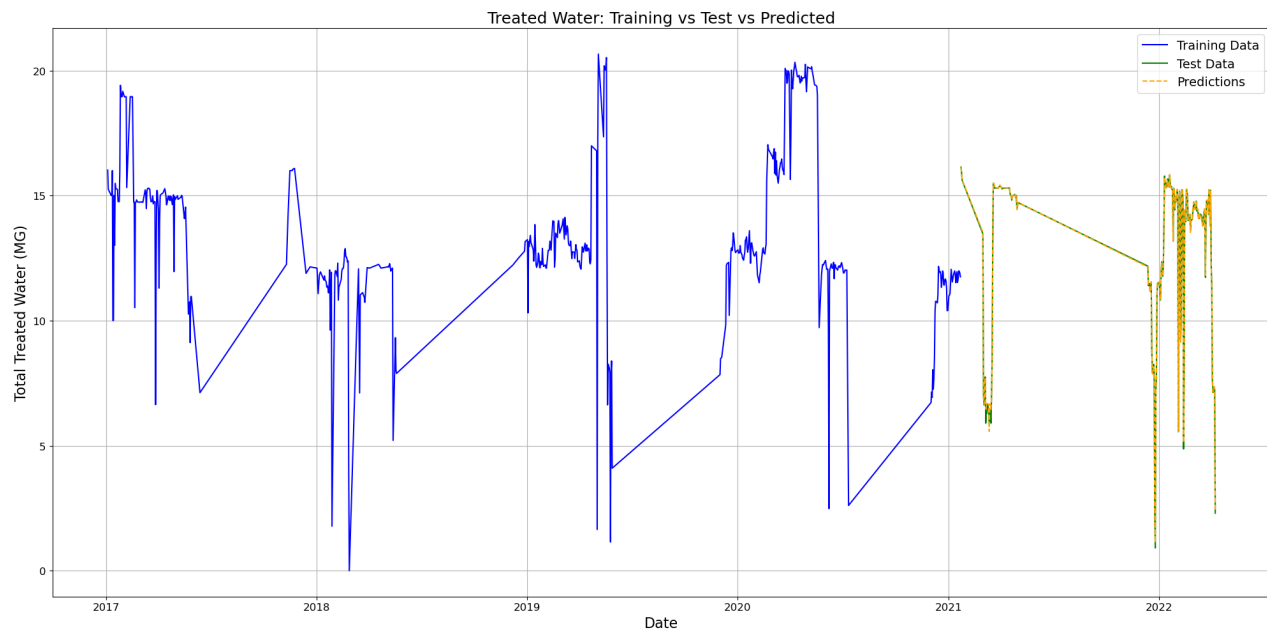
--- Test Set ---
 RMSE: 0.10218589065663168
 MAE: 0.05026823420561527
 Error Rate: 0.01764705882352941

In [130]:

```

1  # Set the figure size
2  plt.figure(figsize=(20, 10))
3
4  # Plot Training Data
5  train_plot, = plt.plot(data['Date'][:len(y_train_treated_actual)], y_train_treated_actual, label='Training Data', color='blue')
6
7  # Plot Test Data
8  test_plot, = plt.plot(data['Date'][len(y_train_treated_actual):len(y_train_treated_actual) + len(y_test_treated_actual)], y_test_trea
9
10 # Overlay Predicted Data
11 predicted_plot, = plt.plot(data['Date'][len(y_train_treated_actual):len(y_train_treated_actual) + len(y_test_treated_actual)], y_pred
12
13 # Add titles, Labels, and Legend with increased font sizes
14 plt.title('Treated Water: Training vs Test vs Predicted', fontsize=18)
15 plt.xlabel('Date', fontsize=16)
16 plt.ylabel('Total Treated Water (MG)', fontsize=16)
17 plt.legend(handles=[train_plot, test_plot, predicted_plot], fontsize=14)
18 plt.grid(True)
19
20 # Increase the size of tick Labels
21 plt.xticks(fontsize=12)
22 plt.yticks(fontsize=12)
23
24 # Show the plot
25 plt.tight_layout()
26 plt.show()
27

```



In [131]:

```
1  # Making predictions on the Training set for Energy Consumption
2  y_pred_train_energy = best_gb_model_energy.predict(X_train_actual)
3
4  # Compute the metrics for the training dataset
5  rmse_train_energy = sqrt(mean_squared_error(y_train_energy_actual, y_pred_train_energy))
6  mae_train_energy = mean_absolute_error(y_train_energy_actual, y_pred_train_energy)
7  error_rate_train_energy = error(y_train_energy_actual, y_pred_train_energy)
8
9  # Print out the metrics for training data
10 print("=== Training Data for Energy Consumption ===")
11 print(f"RMSE: {rmse_train_energy}")
12 print(f"MAE: {mae_train_energy}")
13 print(f"Error Rate: {error_rate_train_energy}")
14
15 # Making predictions on the test set for Energy Consumption
16 y_pred_test_energy = best_gb_model_energy.predict(X_test_actual)
17
18 # Compute the metrics for the test dataset
19 rmse_test_energy = sqrt(mean_squared_error(y_test_energy_actual, y_pred_test_energy))
20 mae_test_energy = mean_absolute_error(y_test_energy_actual, y_pred_test_energy)
21 error_rate_test_energy = error(y_test_energy_actual, y_pred_test_energy)
22
23 # Print out the metrics for test data
24 print("\n=== Test Data for Energy Consumption ===")
25 print(f"RMSE: {rmse_test_energy}")
26 print(f"MAE: {mae_test_energy}")
27 print(f"Error Rate: {error_rate_test_energy}")
28
```

```
=== Training Data for Energy Consumption ===
RMSE: 17.849240272718422
MAE: 5.251089131176049
Error Rate: 0.796078431372549
```

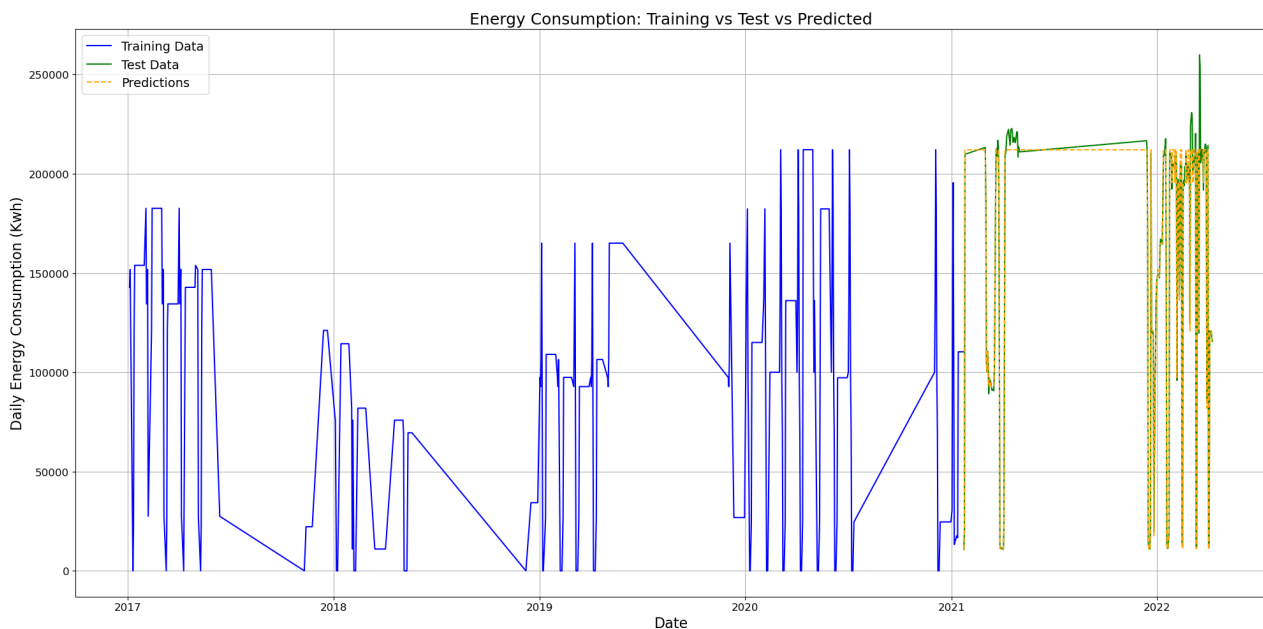
```
=== Test Data for Energy Consumption ===
RMSE: 6751.1928418417865
MAE: 3740.148612606303
Error Rate: 1.0
```

In [133]:

```

1 import matplotlib.pyplot as plt
2
3 # Set the figure size
4 plt.figure(figsize=(20, 10))
5
6 # Plot Training Data for Energy Consumption
7 train_plot, = plt.plot(data['Date'][:len(y_train_energy_actual)], y_train_energy_actual, label='Training Data', color='blue')
8
9 # Plot Test Data for Energy Consumption
10 test_plot, = plt.plot(data['Date'][len(y_train_energy_actual):len(y_train_energy_actual) + len(y_test_energy_actual)], y_test_energy_
11
12 # Overlay Predicted Data for Energy Consumption
13 predicted_plot, = plt.plot(data['Date'][len(y_train_energy_actual):len(y_train_energy_actual) + len(y_test_energy_actual)], y_pred_te
14
15 # Add titles, Labels, and Legend with increased font sizes
16 plt.title('Energy Consumption: Training vs Test vs Predicted', fontsize=18)
17 plt.xlabel('Date', fontsize=16)
18 plt.ylabel('Daily Energy Consumption (Kwh)', fontsize=16)
19 plt.legend(handles=[train_plot, test_plot, predicted_plot], fontsize=14)
20 plt.grid(True)
21
22 # Increase the size of tick labels
23 plt.xticks(fontsize=12)
24 plt.yticks(fontsize=12)
25
26 # Show the plot
27 plt.tight_layout()
28 plt.show()
29

```



In [134]:

```

1 import random
2
3 # Initialize lists to store predictions
4 treated_predictions = []
5 energy_predictions = []

```

In [135]:

```

1 # Number of loops needed to reach 1825 predictions
2 n_loops = 1825 // len(X_train_actual) # Integer division to get full loops
3 remaining_rows = 1825 % len(X_train_actual) # Remaining rows after full loops

```

In [136]:

```

1 # Loop through and concatenate the DataFrame with itself
2 X_future = pd.concat([X_train_actual] * n_loops, ignore_index=True)

```

In [137]:

```

1 # Add the remaining rows
2 if remaining_rows > 0:
3     X_future = pd.concat([X_future, X_train_actual.iloc[:remaining_rows]], ignore_index=True)

```

In [138]:

```
1 # Make the predictions using the GB models
2 pred_treated = best_gb_model_treated.predict(X_future)
3 pred_energy = best_gb_model_energy.predict(X_future)
4
5 # Append the predictions to the lists
6 treated_predictions.extend(pred_treated)
7 energy_predictions.extend(pred_energy)
8
```

In [139]:

```
1 # Convert predictions to DataFrame
2 pred_df = pd.DataFrame({
3     'Treated_Water_Predictions': treated_predictions,
4     'Energy_Consumption_Predictions': energy_predictions
5 })
```

In [140]:

```
1 # Generate a date column
2 last_date = pd.to_datetime(data['Date'].iloc[-1])
3 future_dates = pd.date_range(start=last_date, periods=len(treated_predictions) + 1, freq='D')[1:]
```

In [141]:

```
1 # Add the 'Date' column and rearrange it to be the first column
2 pred_df['Date'] = future_dates
3 pred_df = pred_df[['Date', 'Treated_Water_Predictions', 'Energy_Consumption_Predictions']]
```

In [142]:

```
1 # Export to CSV
2 pred_df.to_csv('Future_Predictions.csv', index=False)
3
4 print("Predictions exported to Future_Predictions.csv")
```

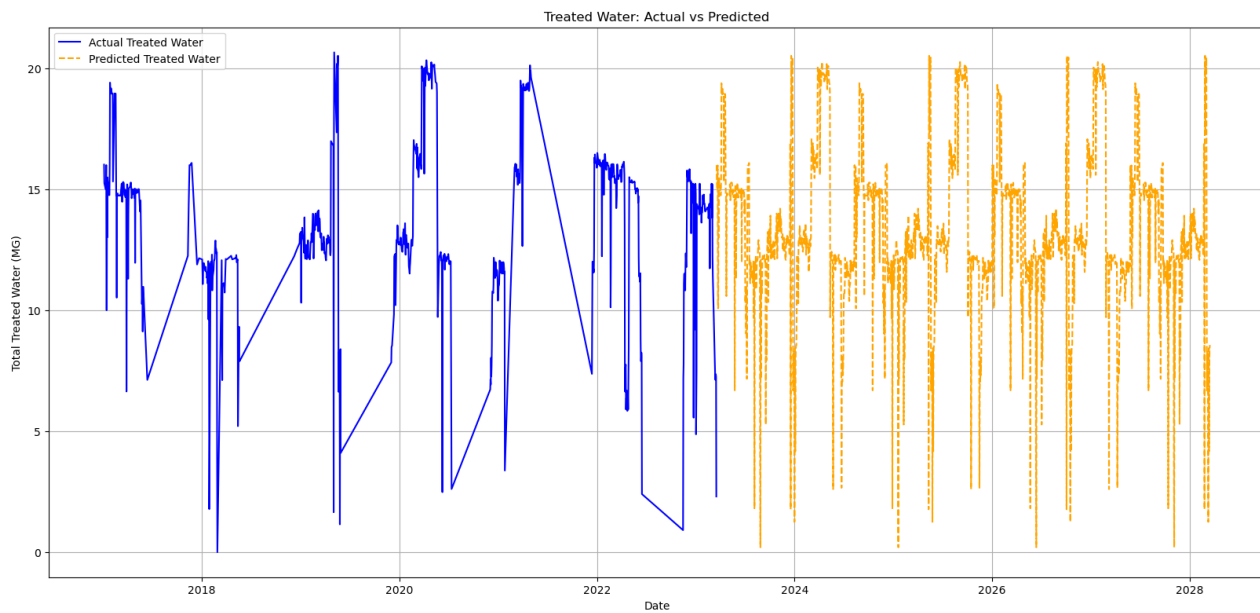
Predictions exported to Future_Predictions.csv

In [144]:

```

1 # Set the figure size
2 plt.figure(figsize=(20, 20))
3
4 # Plot Actual Treated Water
5 plt.subplot(2, 1, 1)
6
7 # Plot the training and test data (actual values)
8 plt.plot(data['Date'], data['TotalTreated water leaving plant,MG'], label='Actual Treated Water', color='blue')
9
10 # Overlay the test data predictions
11 plt.plot(pred_df['Date'], pred_df['Treated_Water_Predictions'], label='Predicted Treated Water', color='orange', linestyle='--')
12
13 plt.title('Treated Water: Actual vs Predicted')
14 plt.xlabel('Date')
15 plt.ylabel('Total Treated Water (MG)')
16 plt.legend()
17 plt.grid(True)
18
19 # Save the figure
20 plt.savefig("Predicted Treated water for Tampa Site.png", dpi=300, bbox_inches='tight')
21
22 plt.show()
23

```

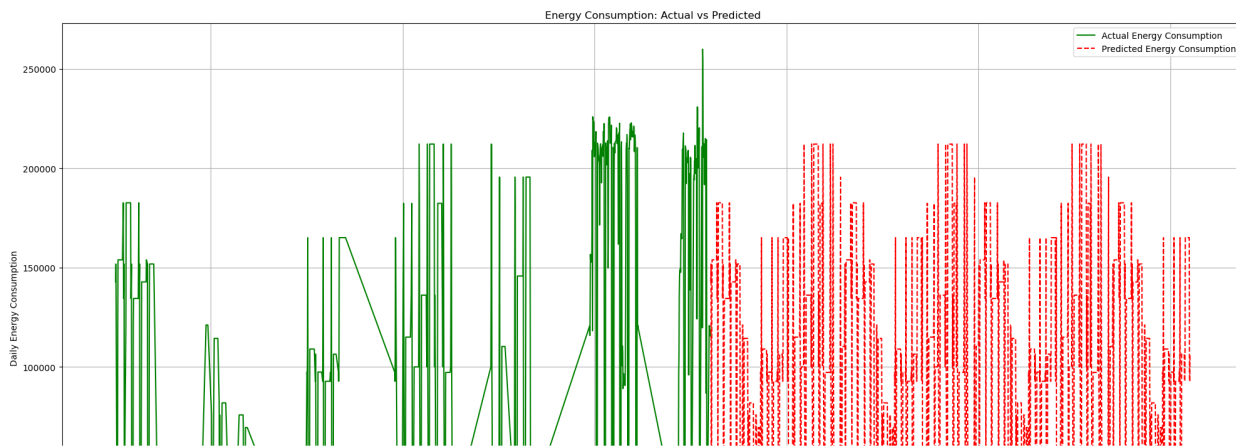


In [145]:

```

1 # Set the figure size
2 plt.figure(figsize=(20, 10))
3
4 # Plot Training + Test Data for Actual Energy Consumption
5 plt.plot(data['Date'], data['Adjusted Daily Consumption, Kwh'], label='Actual Energy Consumption', color='green')
6
7 # Overlay Predicted Energy Consumption for Test Data
8 plt.plot(pred_df['Date'], pred_df['Energy_Consumption_Predictions'], label='Predicted Energy Consumption', color='red', linestyle='--')
9
10 # Add titles, Labels, and Legend
11 plt.title('Energy Consumption: Actual vs Predicted')
12 plt.xlabel('Date')
13 plt.ylabel('Daily Energy Consumption')
14 plt.legend()
15 plt.grid(True)
16
17 # Show the plot
18 plt.tight_layout()
19 plt.show()
20

```



In [146]:

```

1 import pickle
2
3 # Saving the Gradient Boosting model for treated water in Tampa
4 with open('gb_model_treated_tampa.pkl', 'wb') as f:
5     pickle.dump(best_gb_model_treated, f)
6
7 # Saving the Gradient Boosting model for energy consumption
8 with open('best_gb_model_energy_tampa.pkl', 'wb') as f:
9     pickle.dump(best_gb_model_energy, f)
10

```

In []:

```

1 # Loading the Gradient Boosting model for treated water
2 with open('gb_model_treated_tampa.pkl', 'rb') as f:
3     loaded_gb_model_treated_tampa = pickle.load(f)
4
5 # Loading the Gradient Boosting model for energy consumption
6 with open('gb_model_energy_tampa.pkl', 'rb') as f:
7     loaded_gb_model_energy_tampa = pickle.load(f)
8

```