

PJM HOURLY ENERGY CONSUMPTION

PROJECT - 01 : GROUP - 03

FINAL REVIEW PRESENTATION

GROUP MEMBERS -

DEEPESH, DEVA DHARSHINI, SUJAL,
SHIVA, VIVEK, SRAVYA.



BUSINESS OBJECTIVE :

ABOUT DATASET :-

PJM Interconnection LLC (PJM) is a regional transmission organization (RTO) in the United States. It is part of the Eastern Interconnection grid operating an electric transmission system serving all or parts of Delaware, Illinois, Indiana, Kentucky, Maryland, Michigan, New Jersey, North Carolina, Ohio, Pennsylvania, Tennessee, Virginia, West Virginia, and the District of Columbia. The hourly power consumption data comes from PJM's website and are in megawatts (MW).

PROBLEM STATEMENT :-

1. Split the last year into a test set- can you build a model to predict energy consumption.
2. Find trends in energy consumption around hours of the day, holidays, or long term trends.
3. Understand how daily trends change depending of the time of year. Summer trends are very different than winter trends.
4. Forecast for next 30 days.

PROJECT FLOW :

1. Understanding the problem statement
2. EDA (Exploratory Data Analysis)
 - Data Pre-processing
 - Data cleaning
 - Visualization
 - Weekday, holiday and seasonal Trend
 - Outlier Detection
 - Outlier Removal
1. Model Building
 - 5 Models
 - Choosing The Final Model
 - Forecast For Next 30 Days
1. Model Deployment

PHASE - 1

EXPLORATORY DATA ANALYSIS (EDA)

Data Information

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 143206 entries, 0 to 143205
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Datetime    143206 non-null   datetime64[ns]
 1   PJMW_MW    143206 non-null   float64 
dtypes: datetime64[ns](1), float64(1)
memory usage: 2.2 MB
```

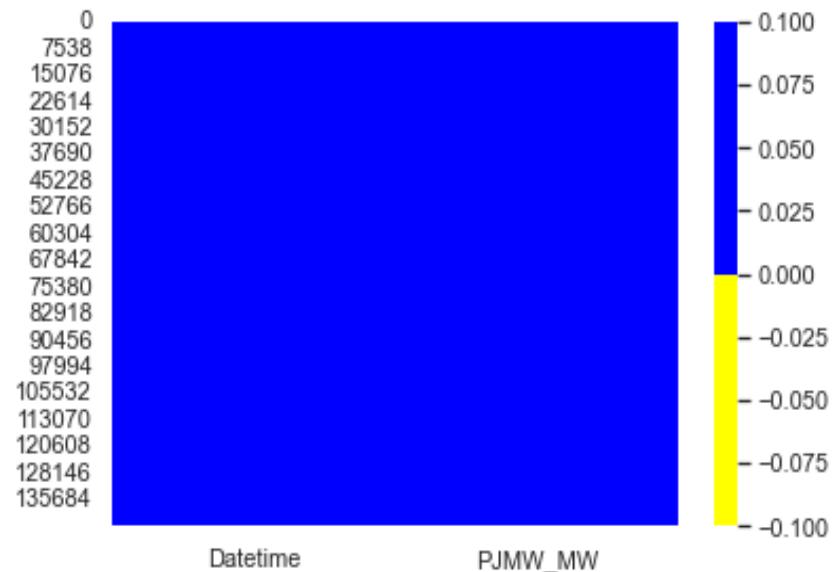
Data Description

```
In [9]: df.describe()
```

```
Out[9]:
```

PJMW_MW	
count	143206.000000
mean	5602.375089
std	979.142872
min	487.000000
25%	4907.000000
50%	5530.000000
75%	6252.000000
max	9594.000000

Checking for Null Values



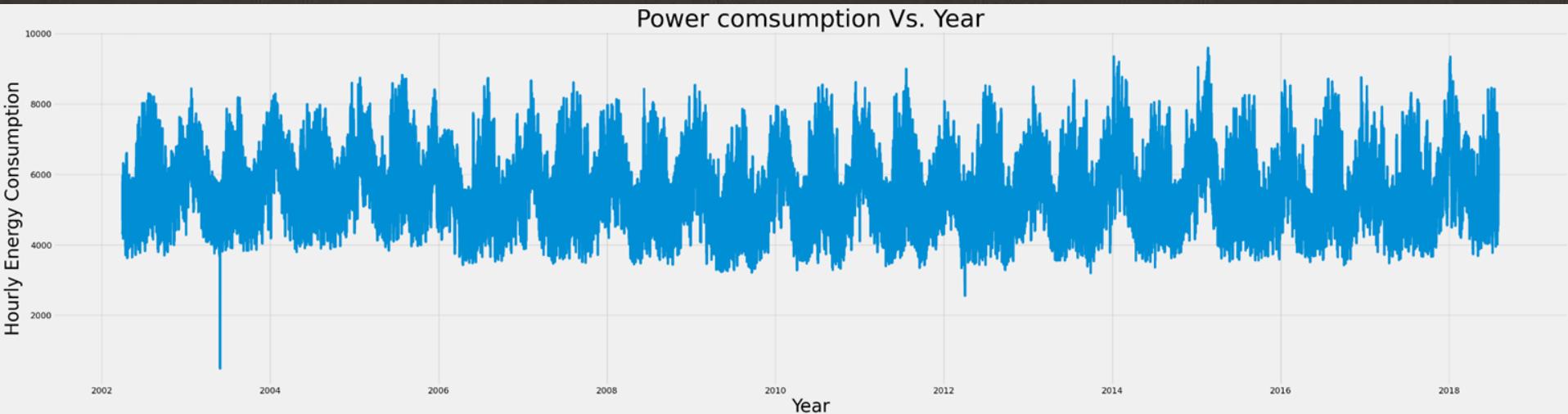
```
In [4]: df.duplicated().sum()
```

```
Out[4]: 0
```

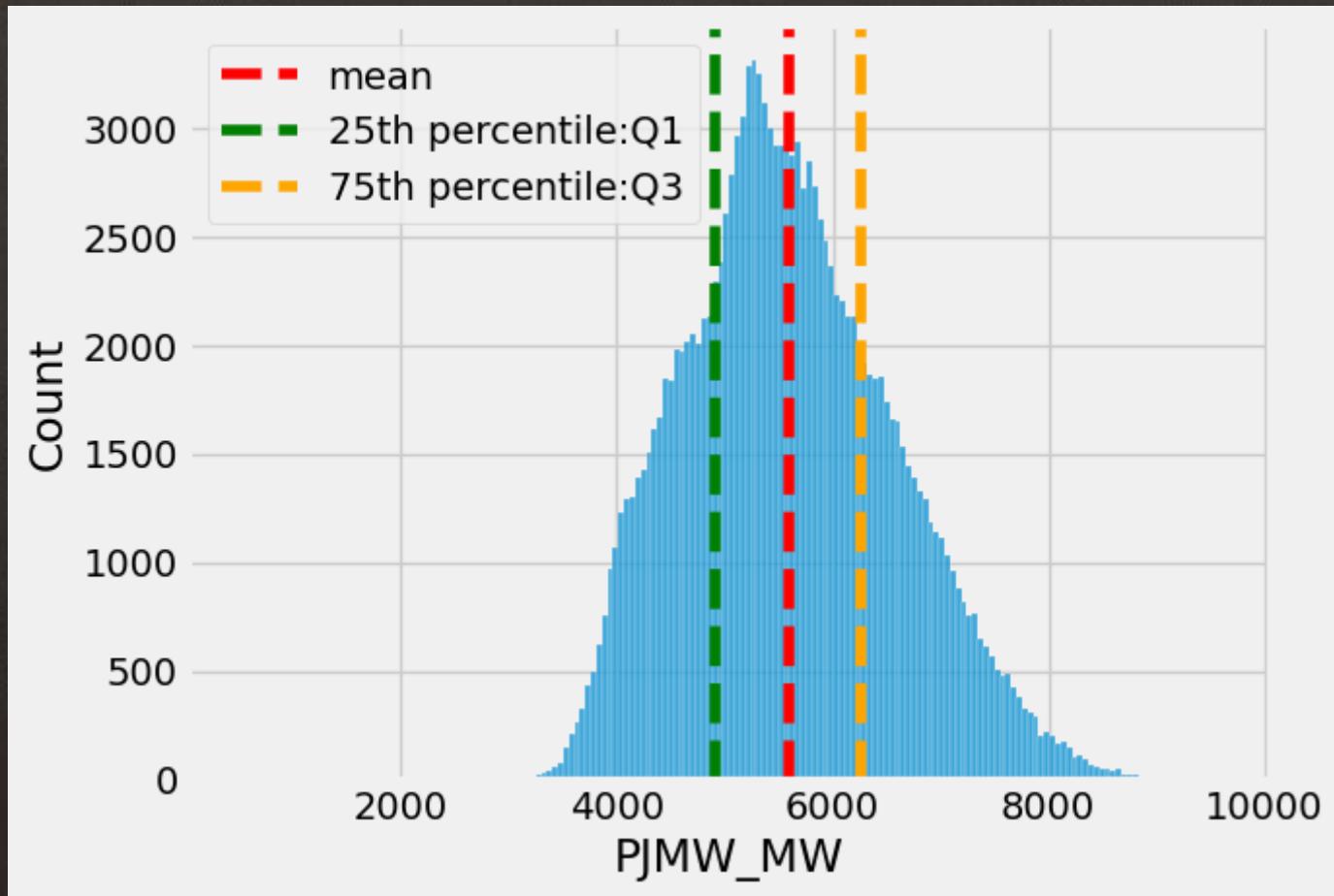
```
In [11]: df.isnull().sum()
```

```
Out[11]: Datetime      0  
PJMW_MW        0  
dtype: int64
```

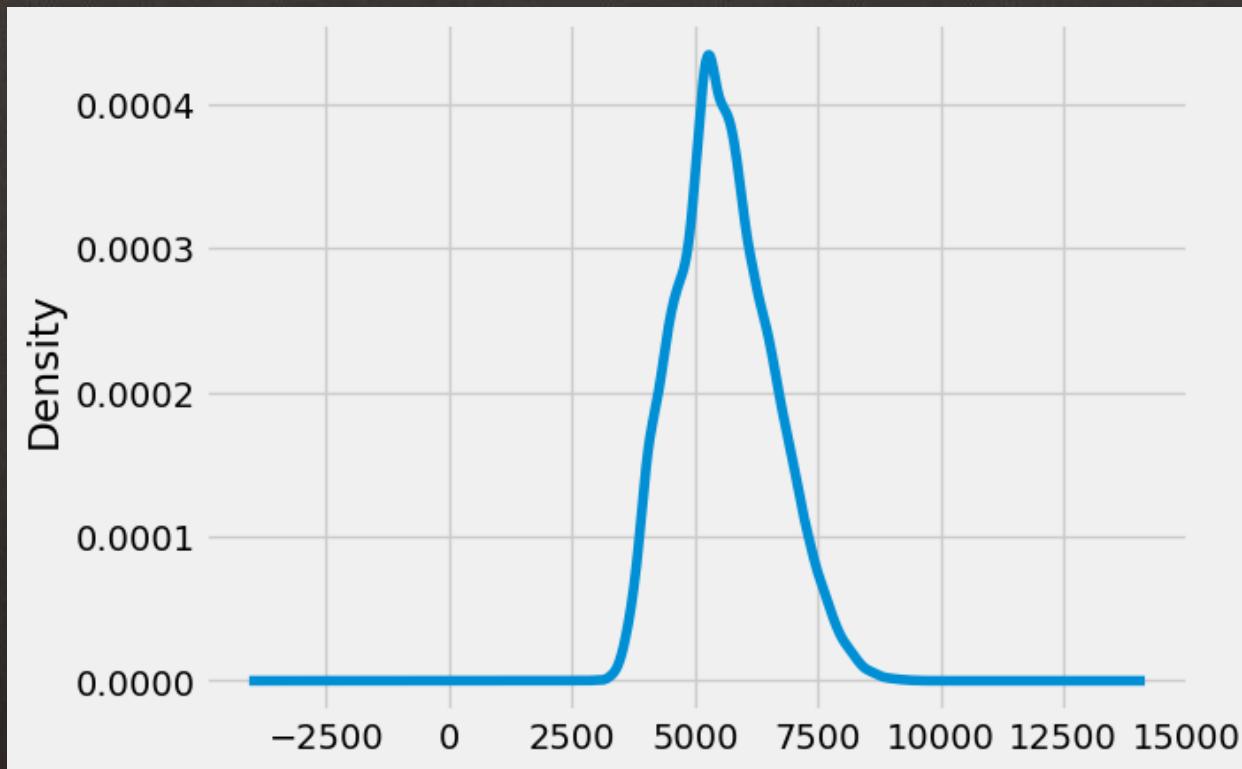
Plotting the data



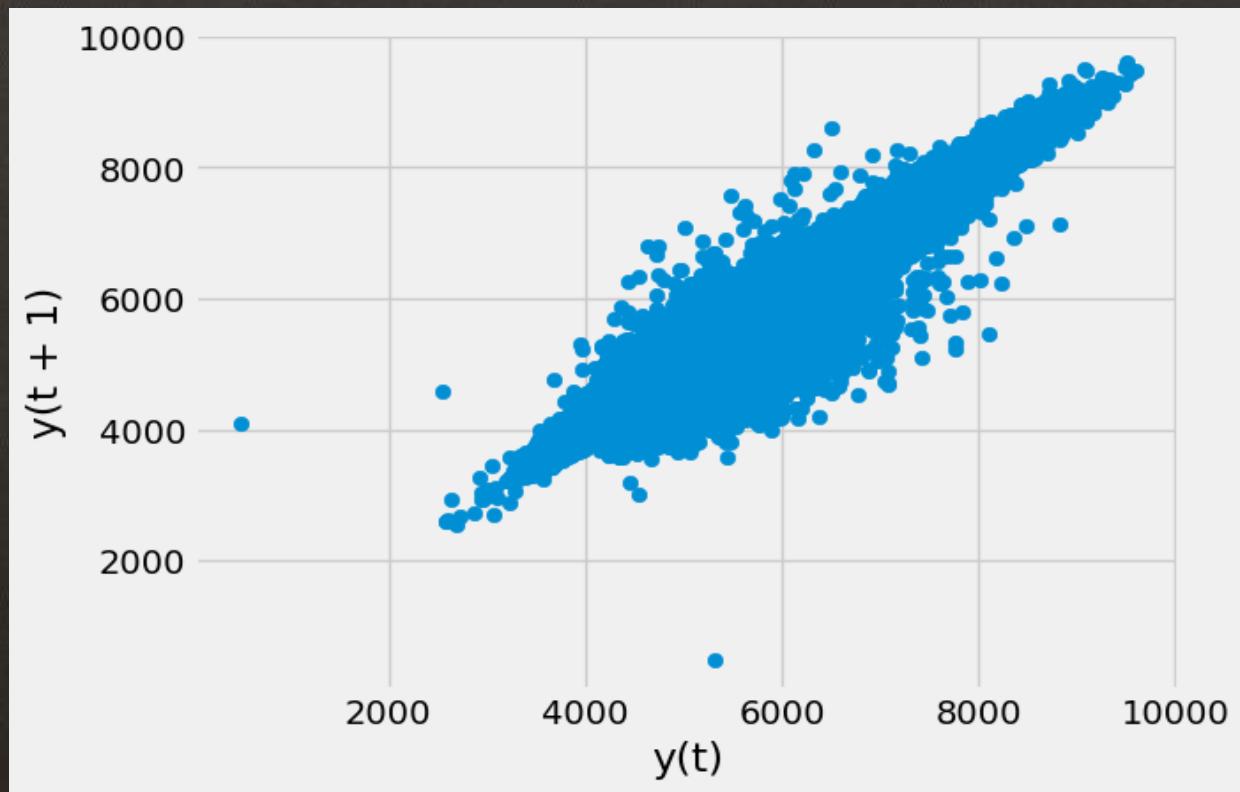
Histogram



Density Plot



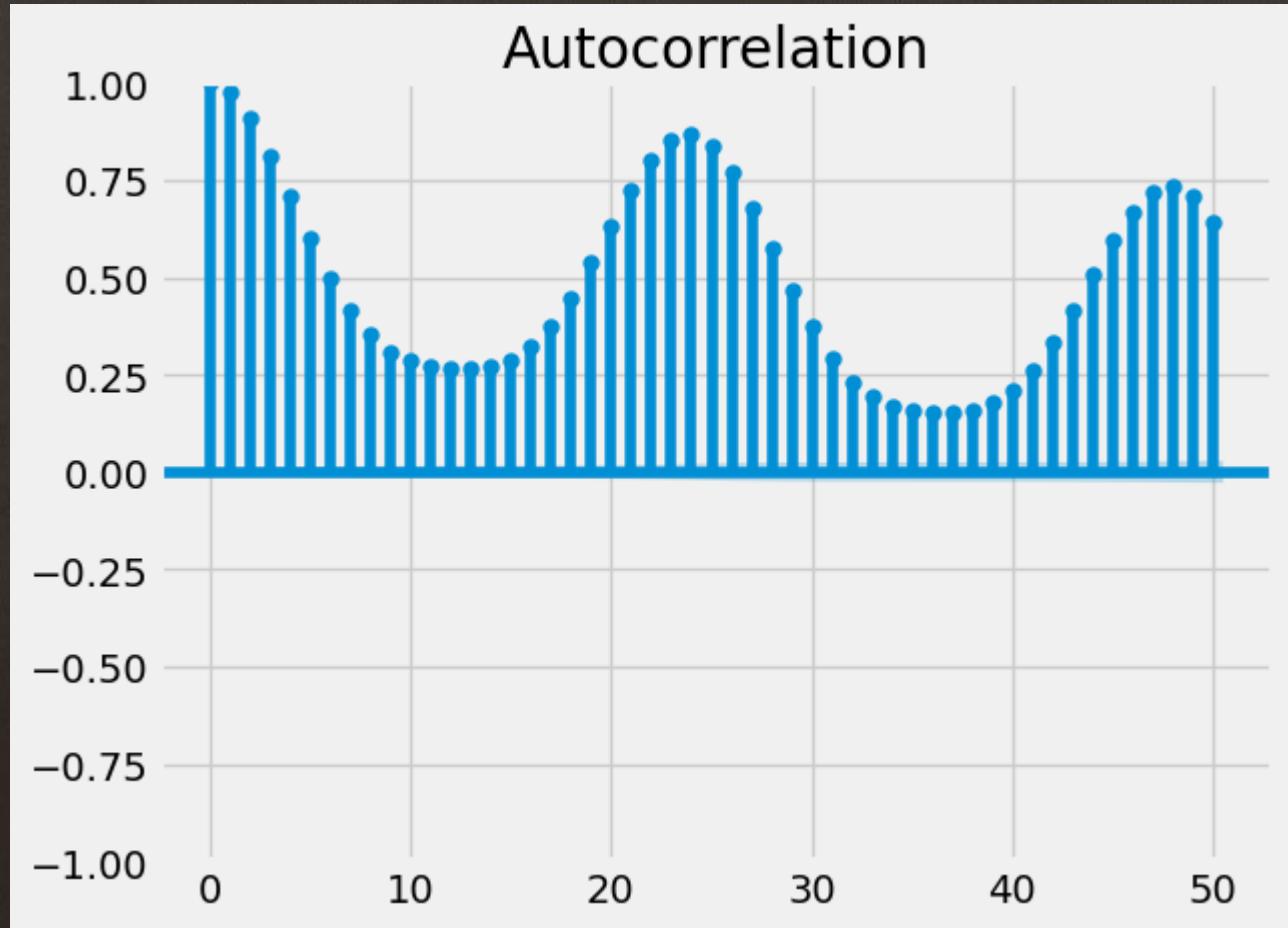
Lag Plot



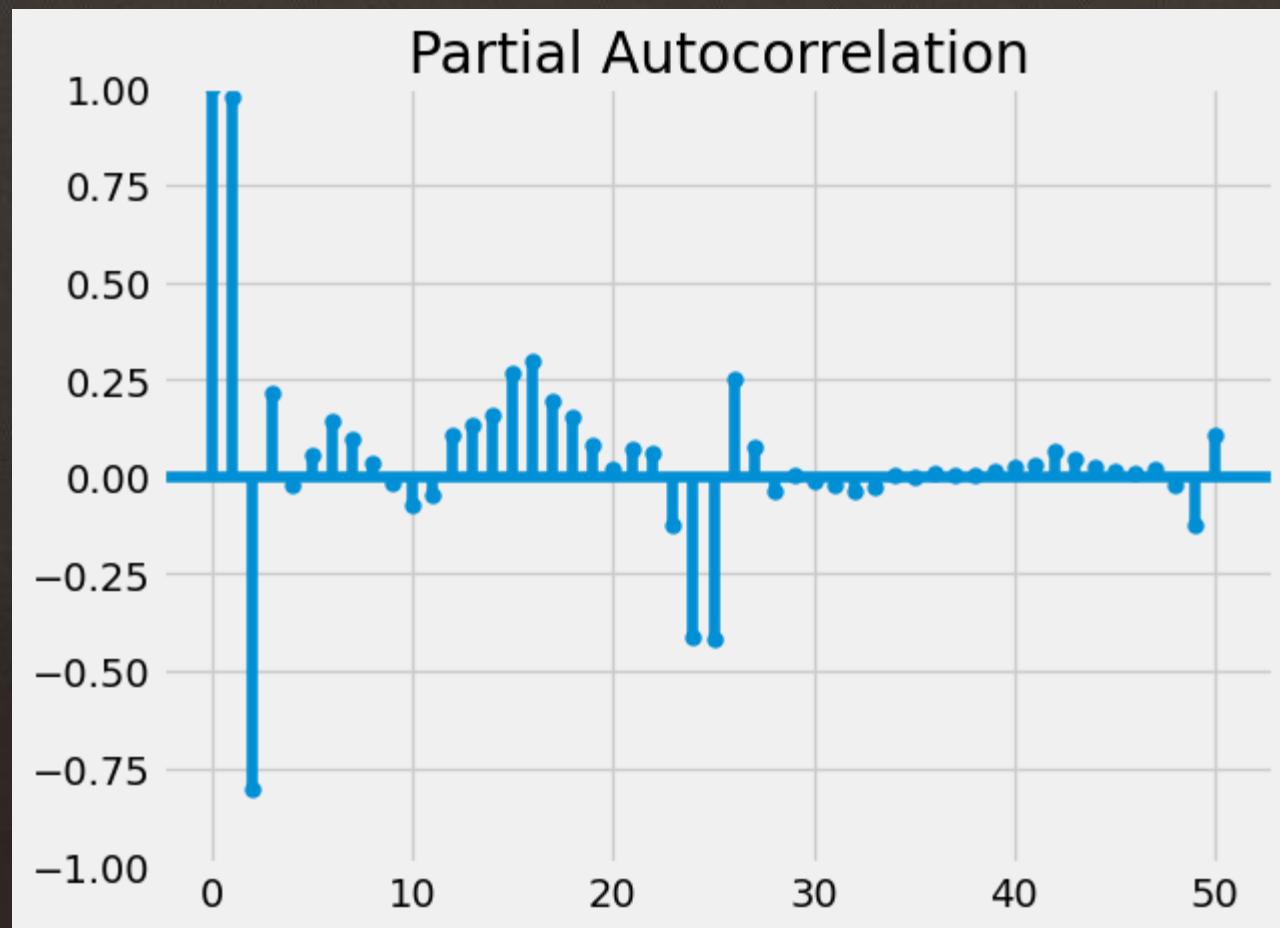
Heatmap



ACF Plot

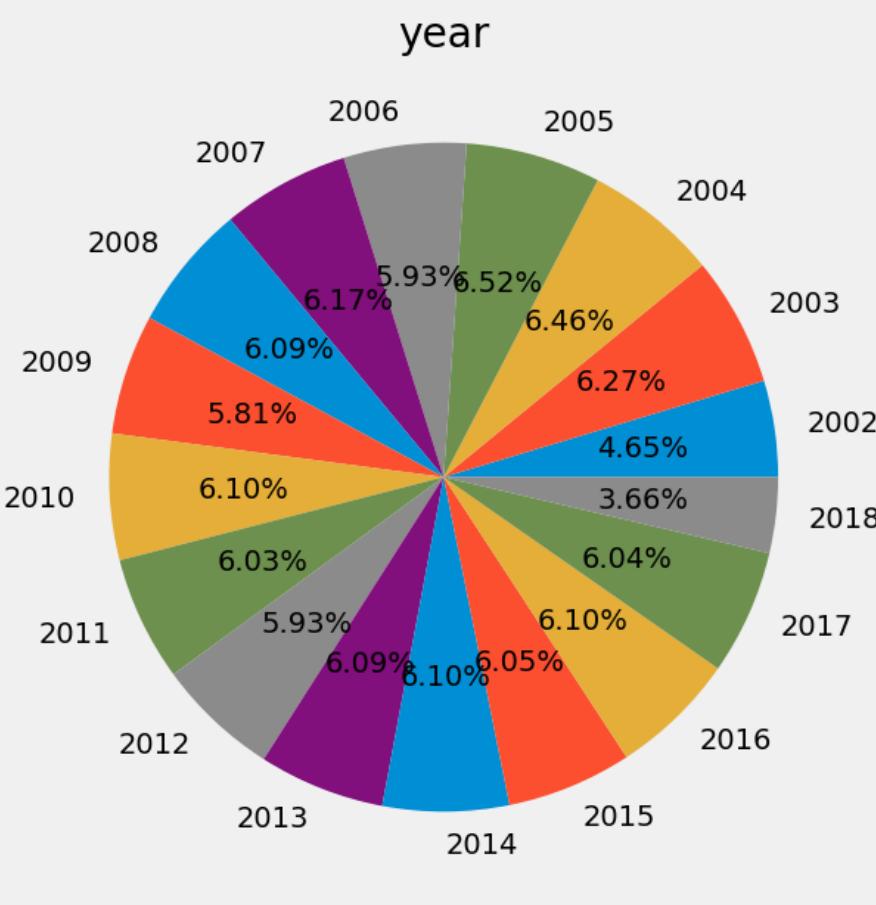


PACF Plot

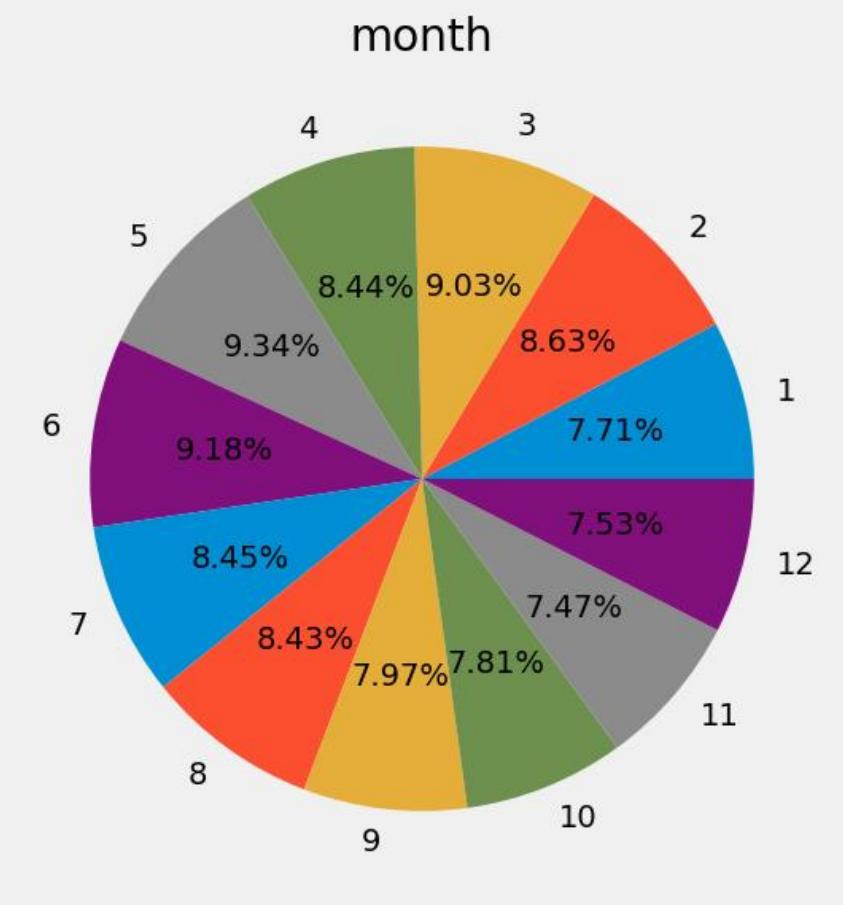


Pie Chart

year

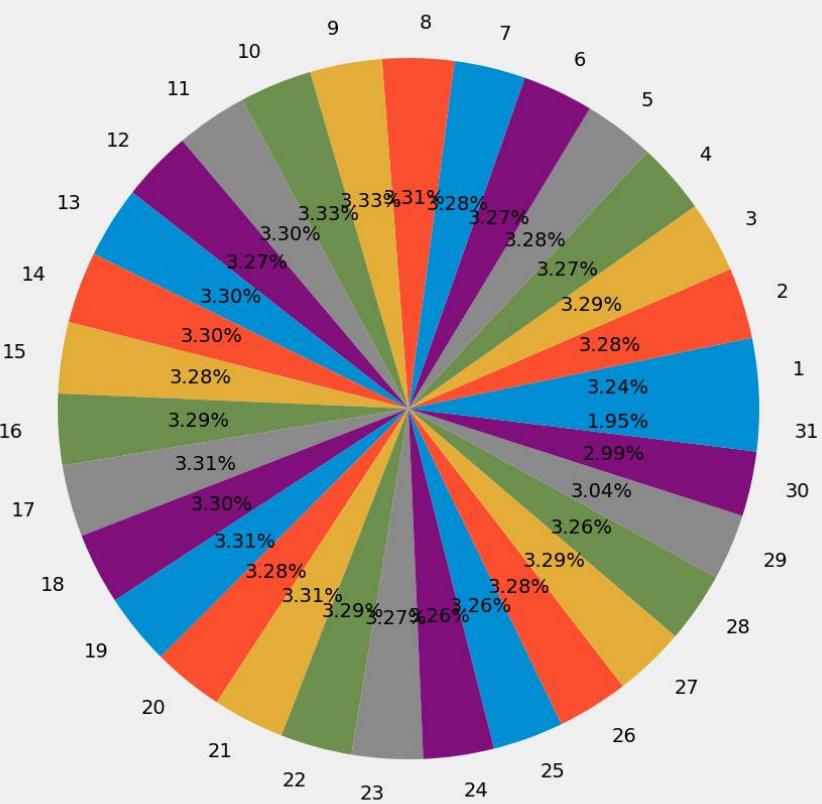


month

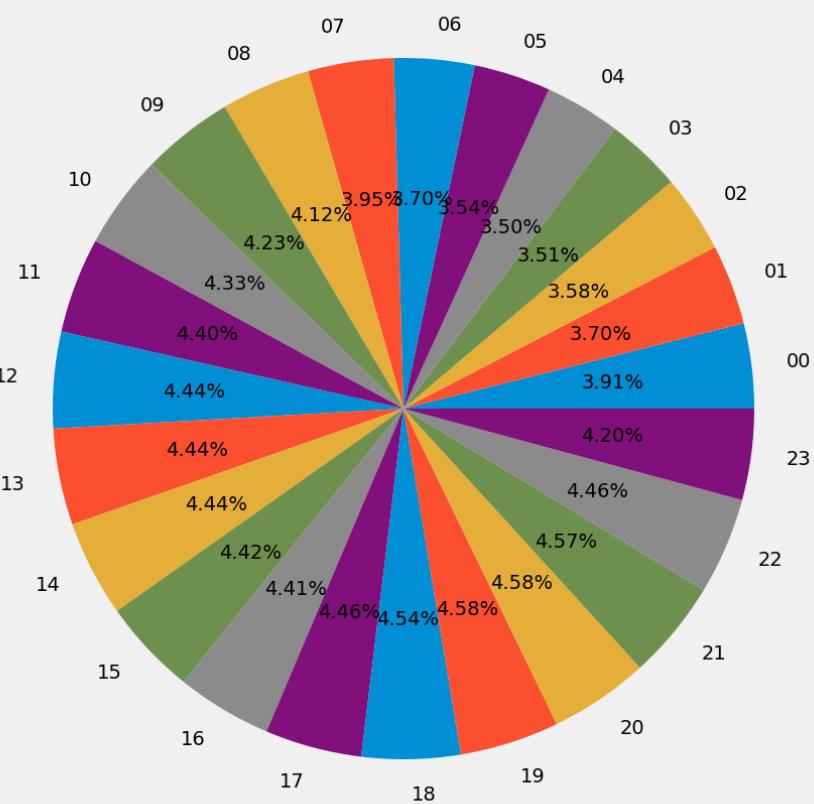


Pie Chart

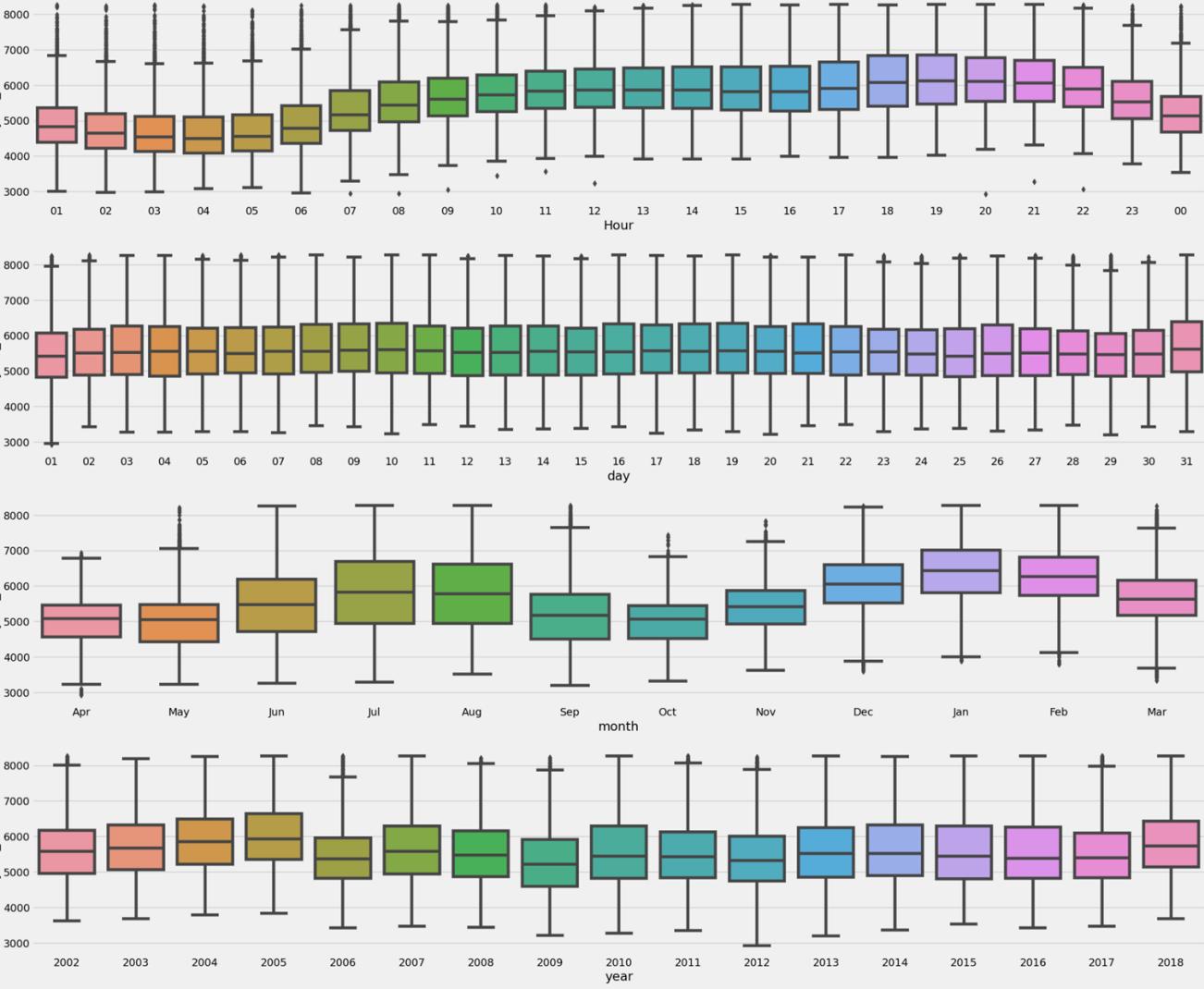
day



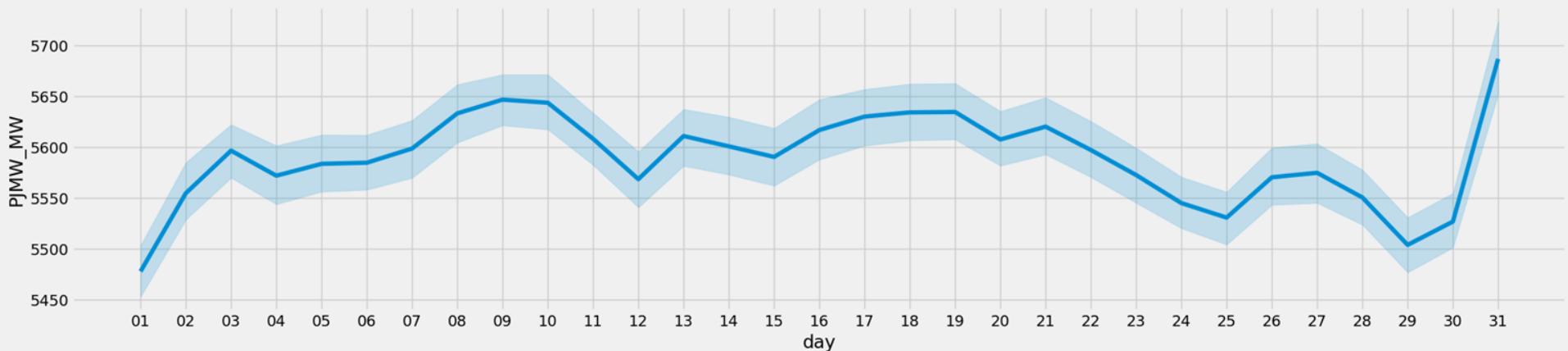
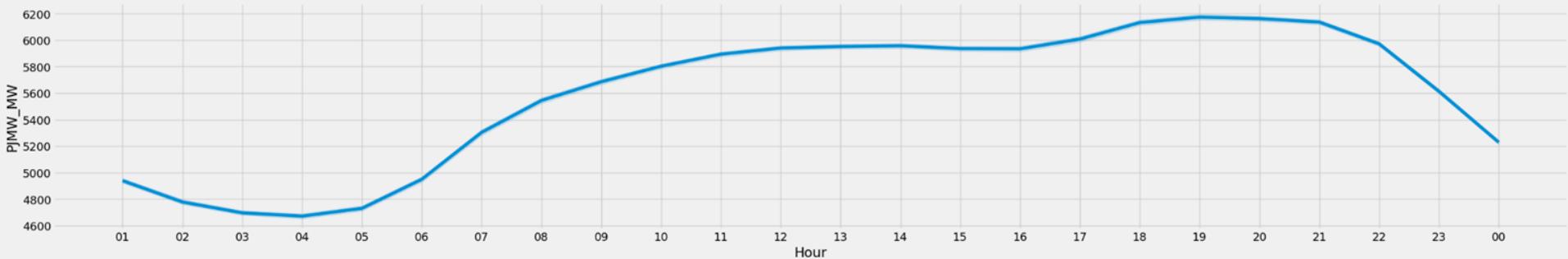
hour

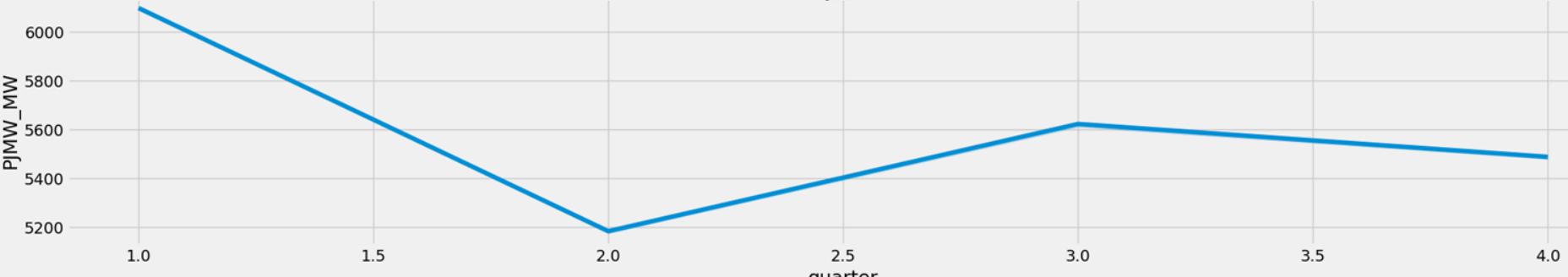
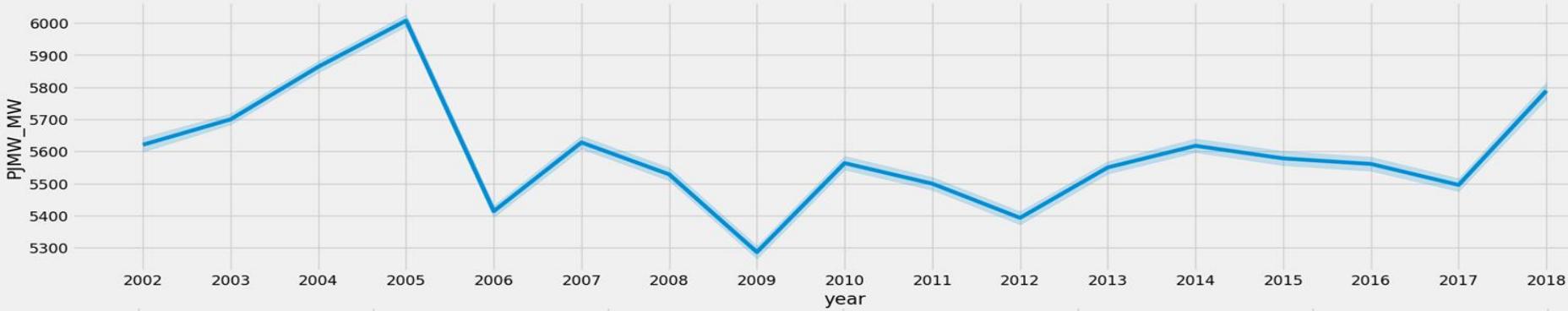


Box Plot



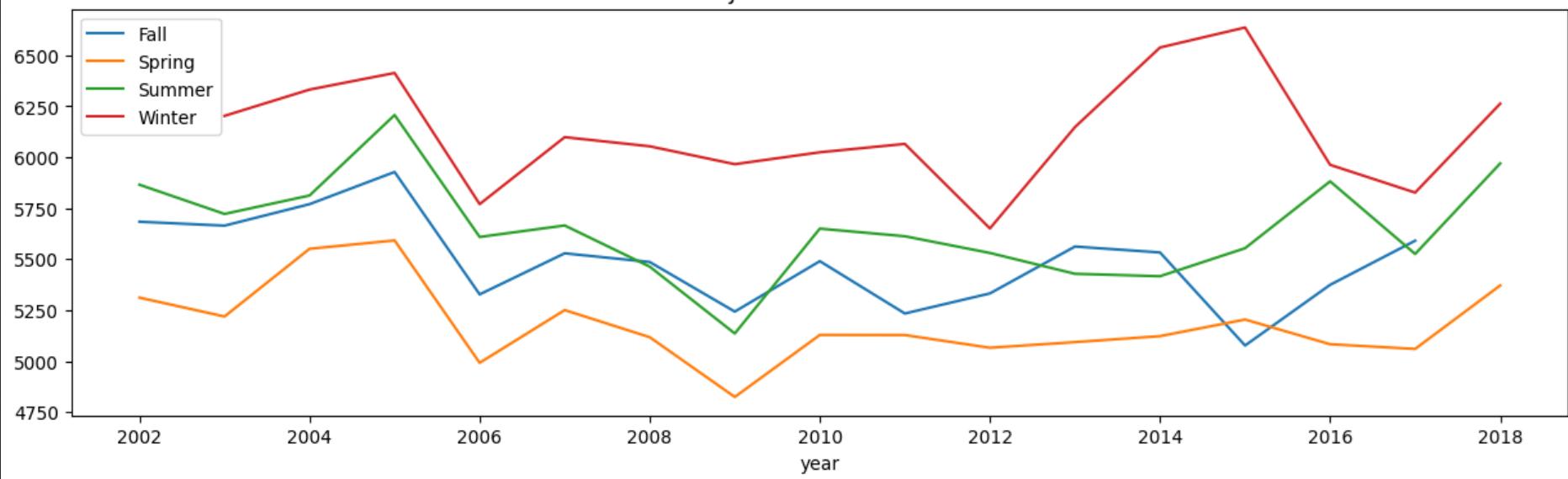
Line Plot



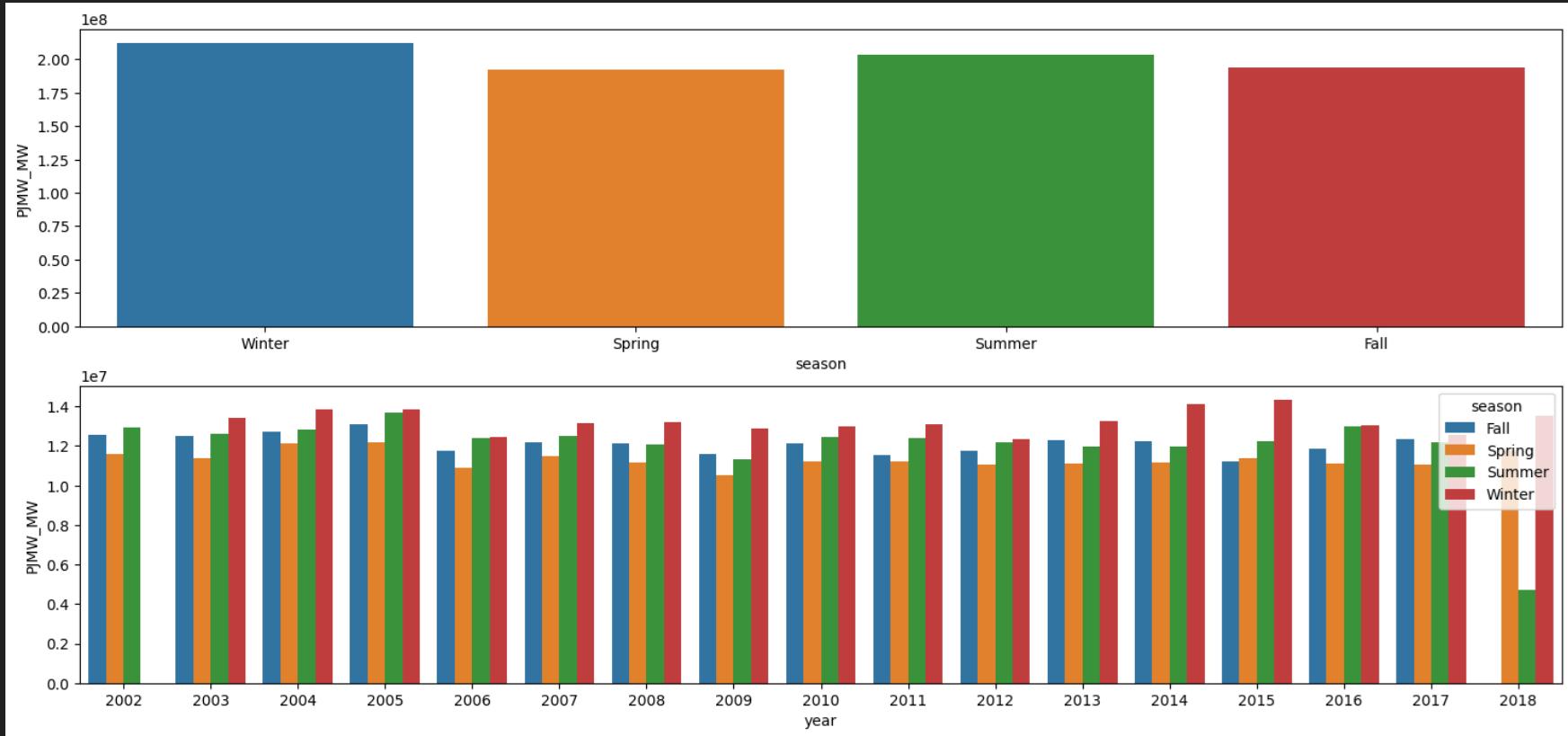


Seasonal Plot

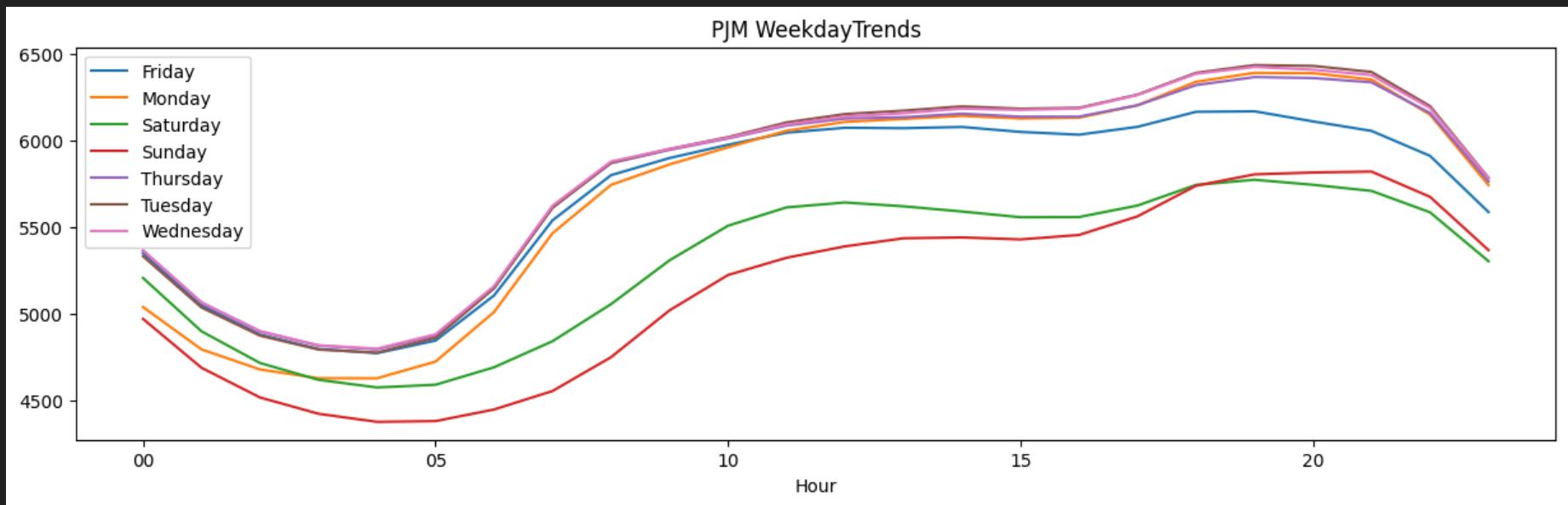
PJM Seasonal Trend



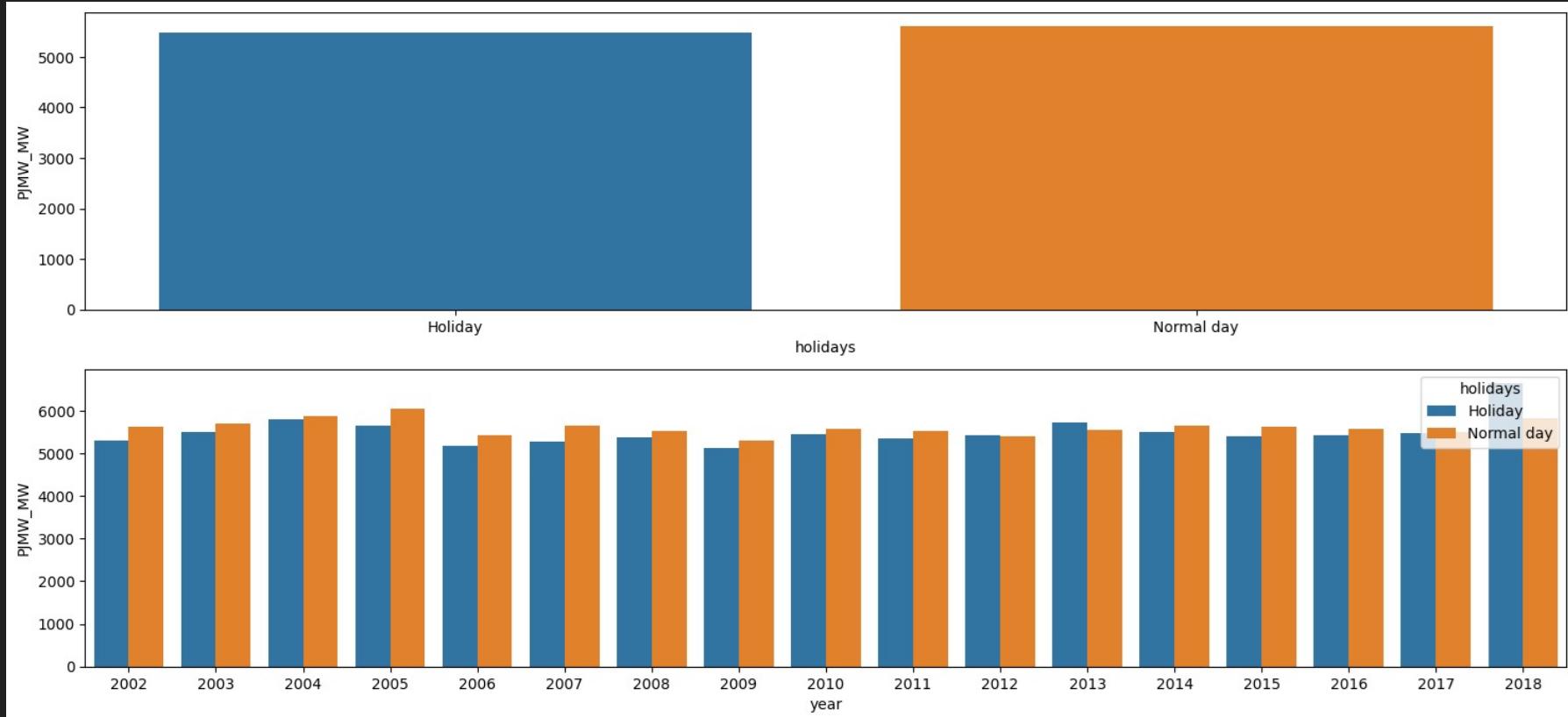
Seasonal Trend



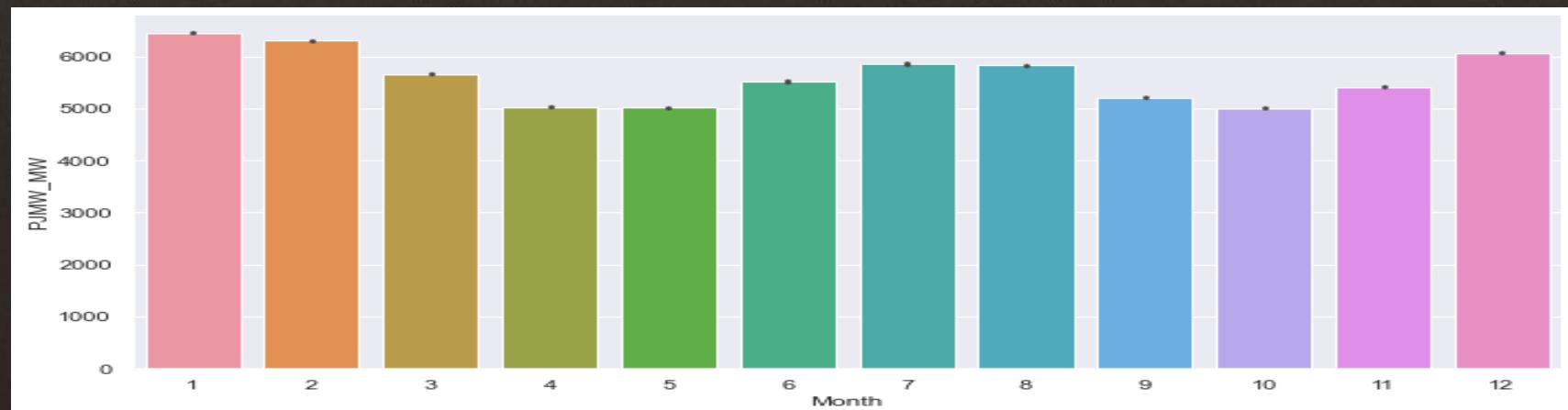
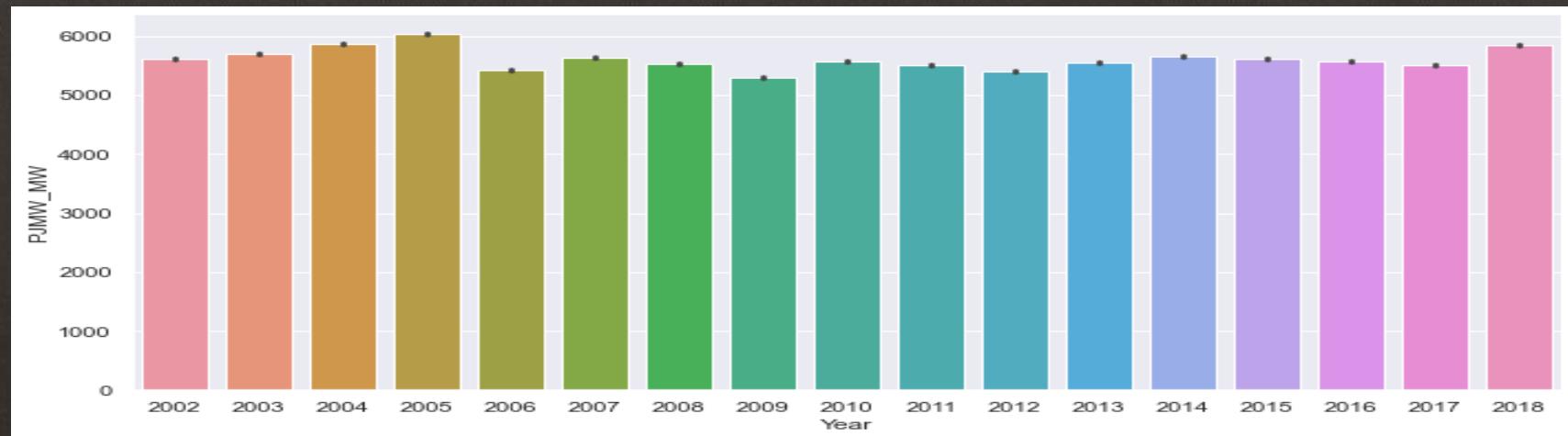
Weekday Trends

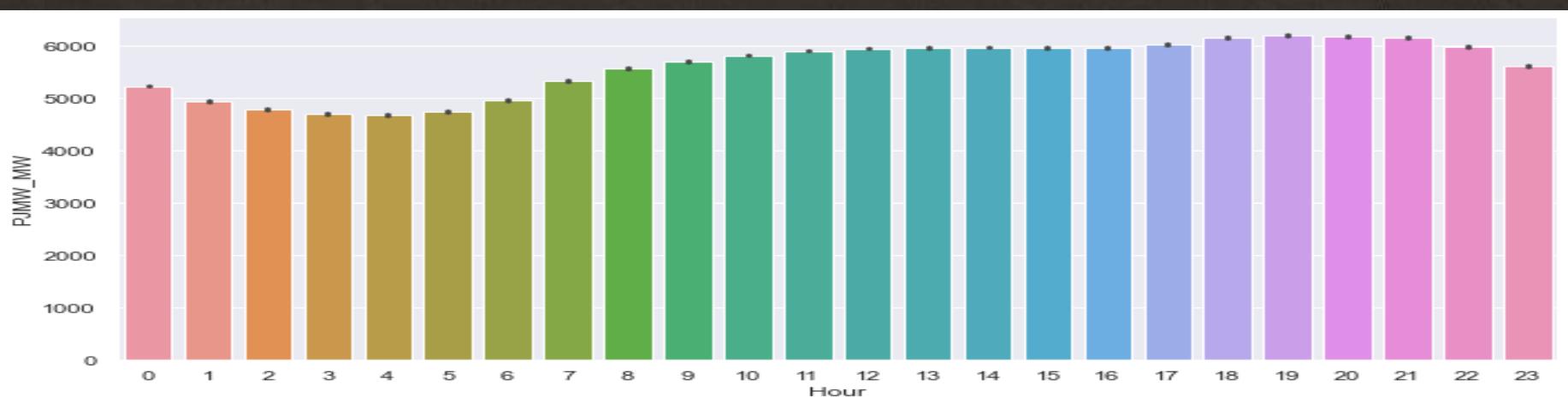
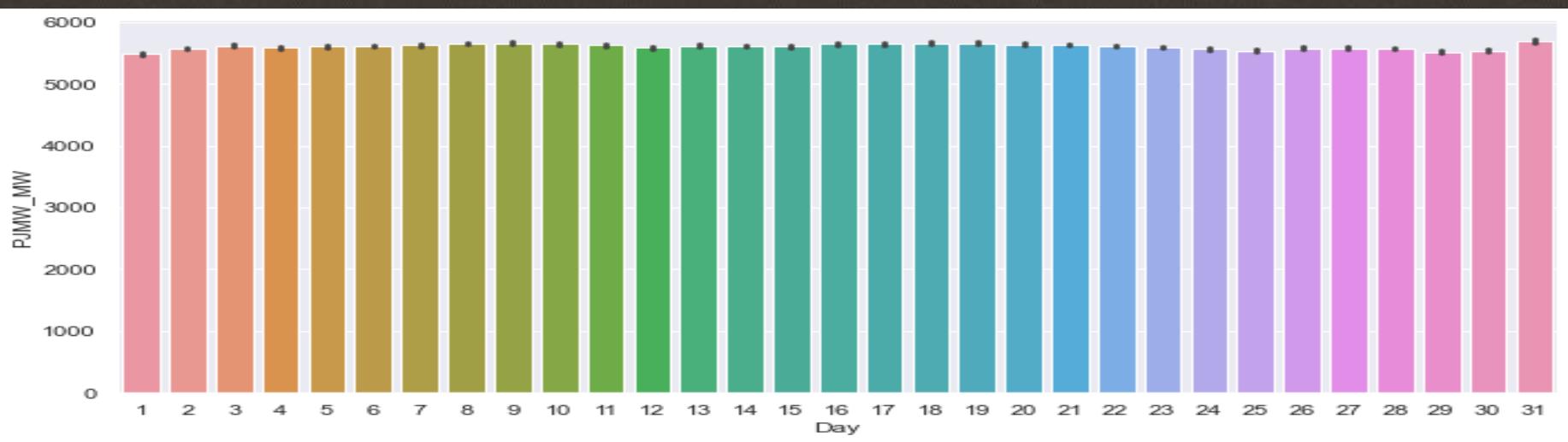


Trend Of Energy Consumption During Holidays

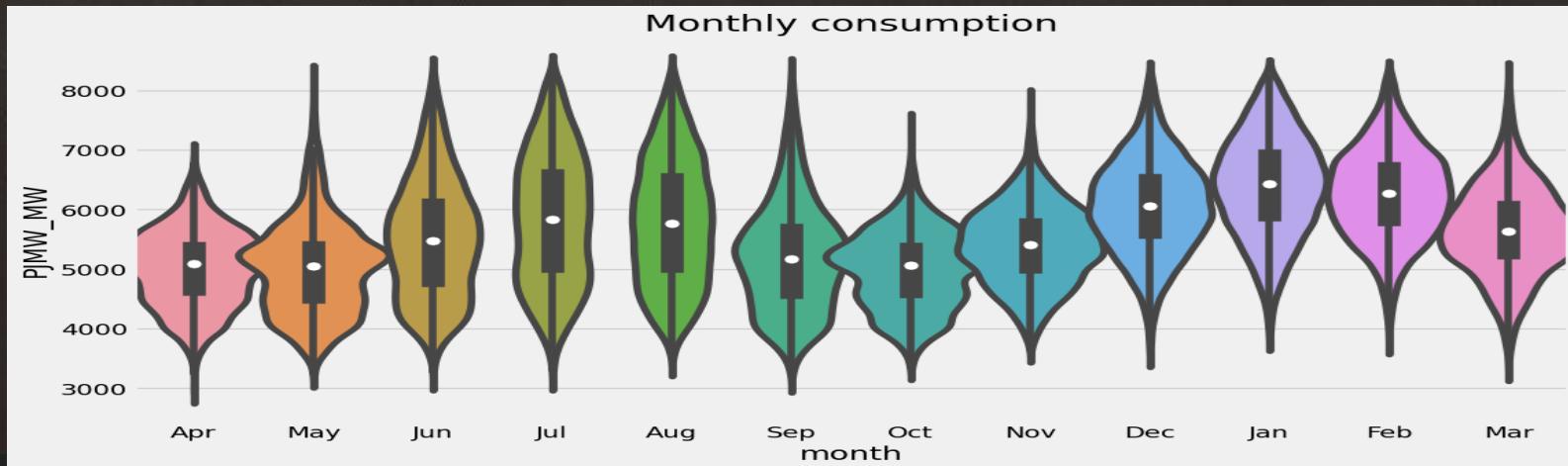
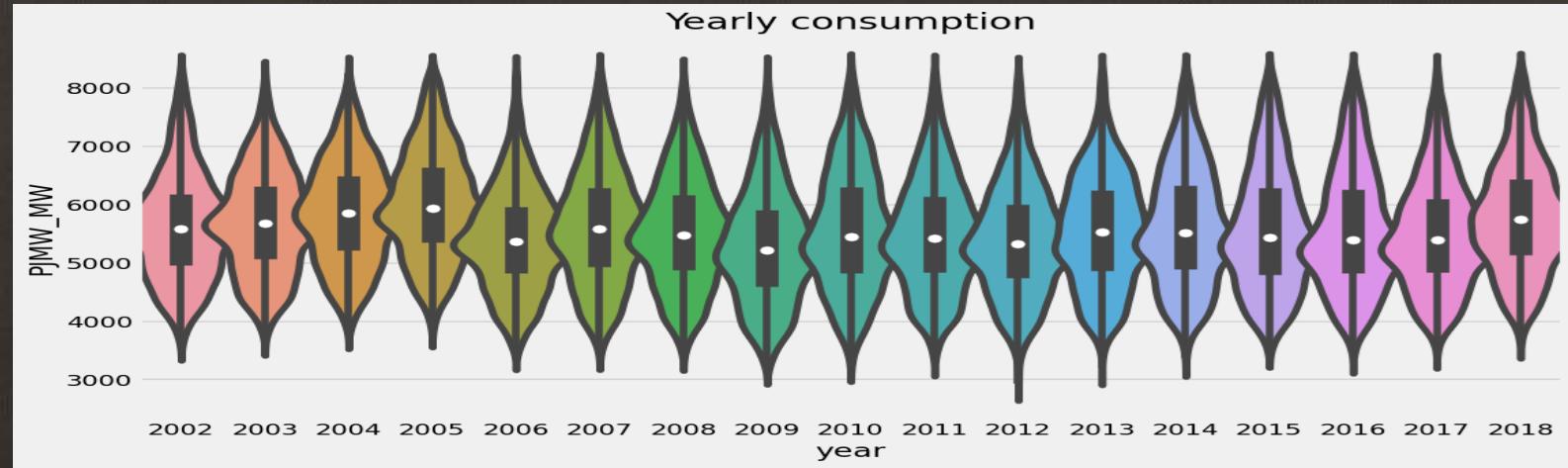


Cat Plots

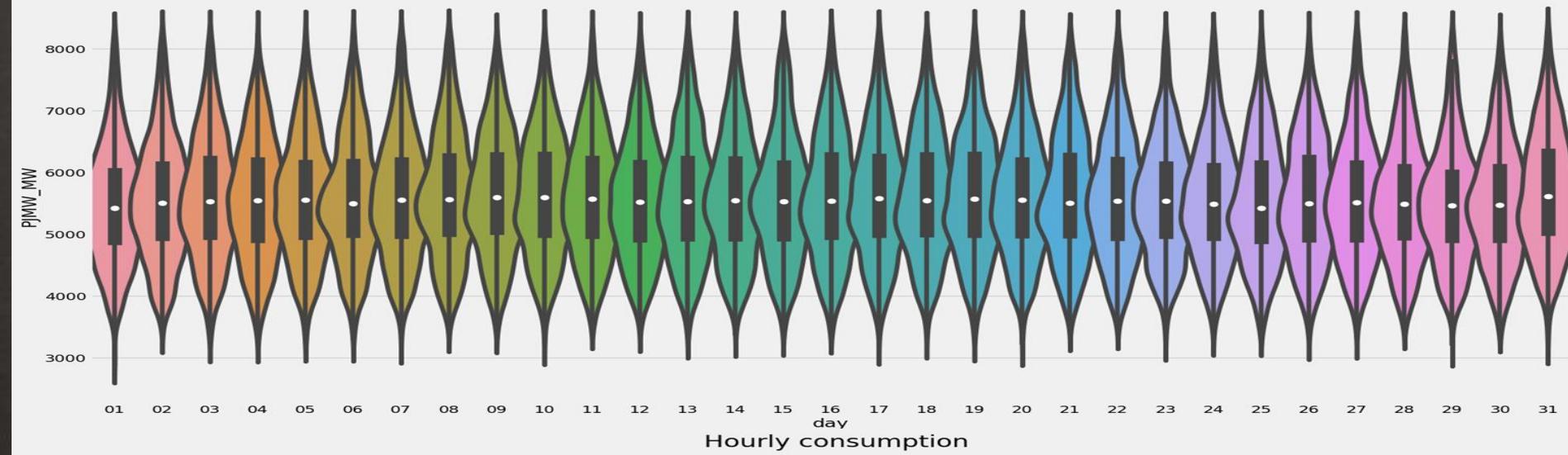




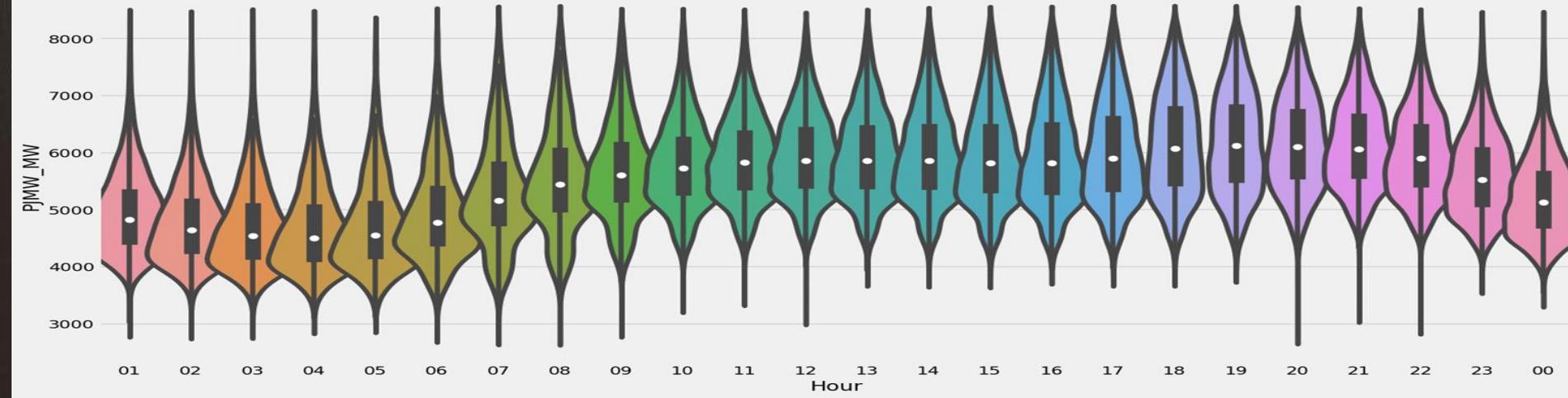
Violin Plot



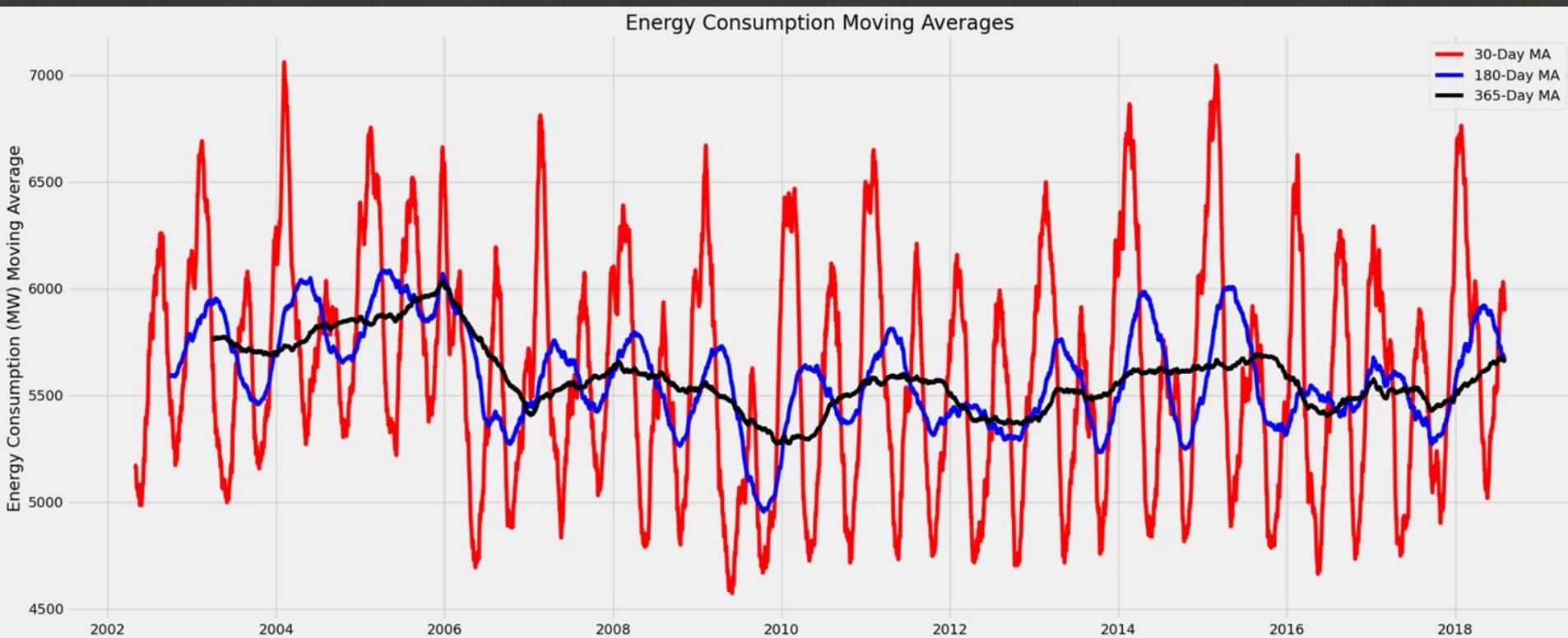
Daily consumption



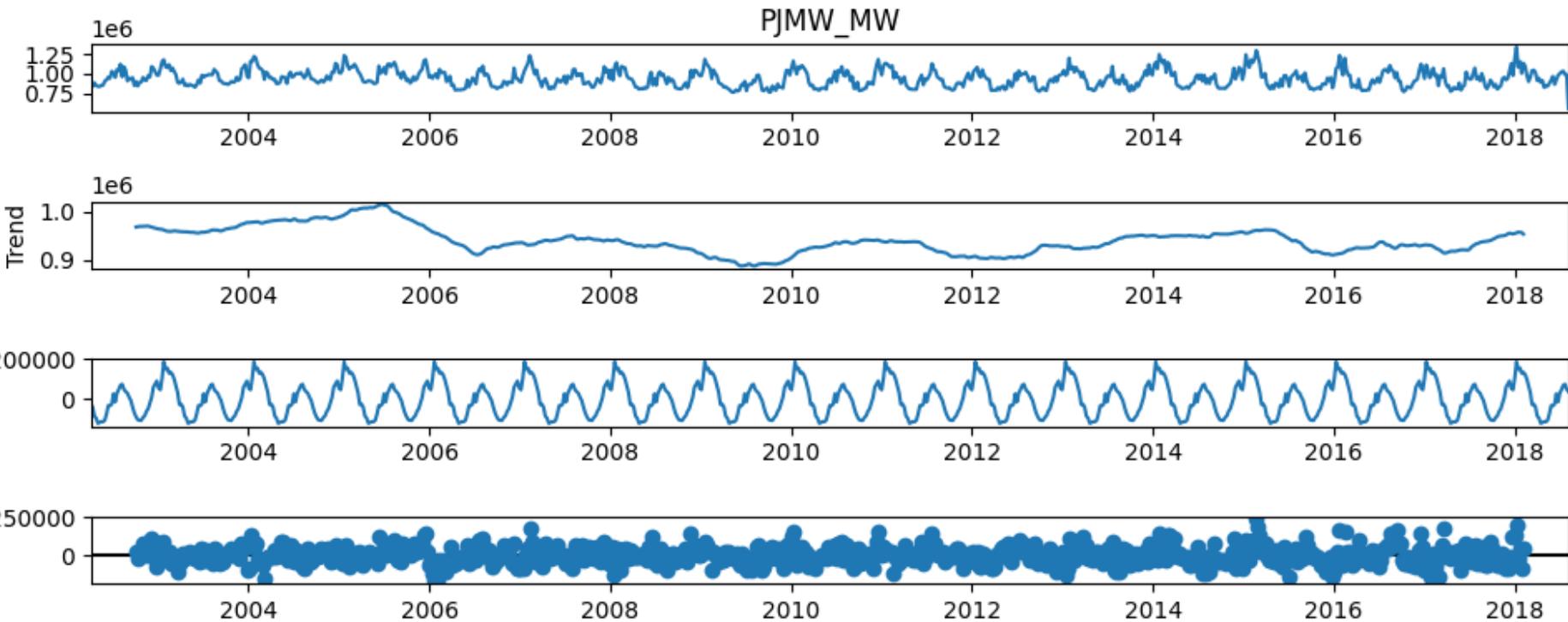
Hourly consumption



Moving Average

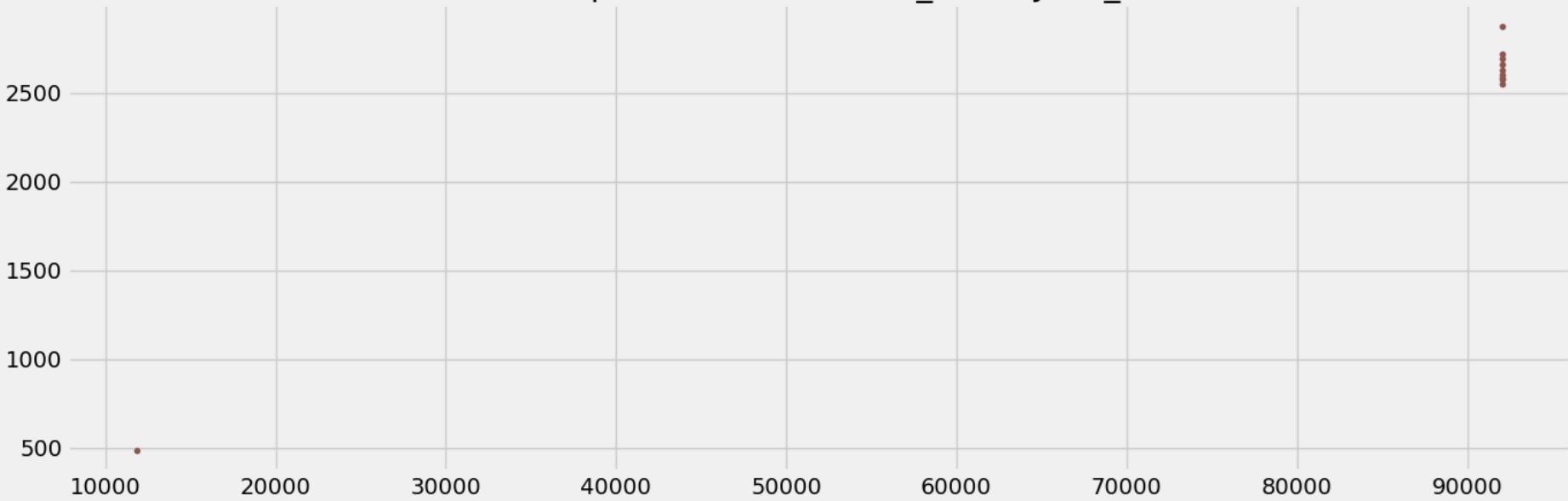


Time Series Decomposition

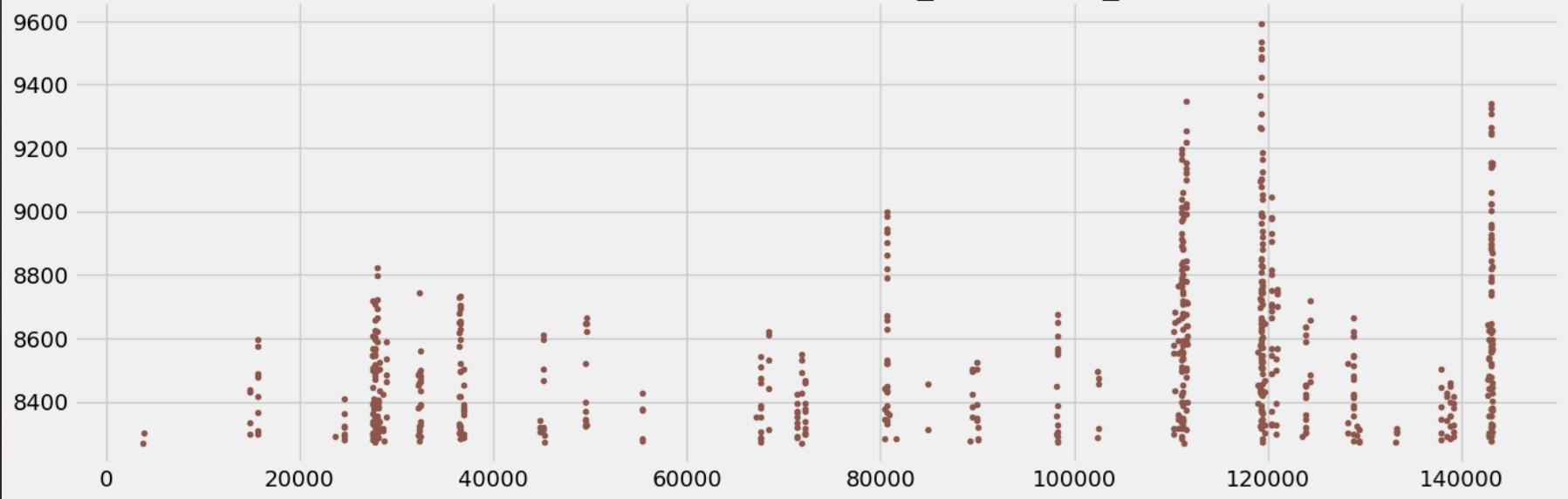


Outlier Detection

Outliers which is present below lower_limit(PJMW_MW- 2889.5)



Outliers which is present above upper_limit(PJMW_MW- 8269.5)



Outlier Removal

```
In [28]: Q1 = df.PJMW_MW.quantile(0.25)
          Q3 = df.PJMW_MW.quantile(0.75)
          Q1, Q3
```

```
Out[28]: (4907.0, 6252.0)
```

```
In [29]: IQR = Q3 - Q1
          IQR
```

```
Out[29]: 1345.0
```

```
In [30]: lower_limit = Q1 - 1.5*IQR
          upper_limit = Q3 + 1.5*IQR
          lower_limit, upper_limit
```

```
Out[30]: (2889.5, 8269.5)
```

```
In [31]: df[(df.PJMW_MW < lower_limit) | (df.PJMW_MW > upper_limit)]
          df.shape
```

```
Out[31]: (143206, 1)
```

```
In [32]: df = df[(df.PJMW_MW > lower_limit) & (df.PJMW_MW < upper_limit)]
          df.shape
```

```
Out[32]: (142503, 1)
```

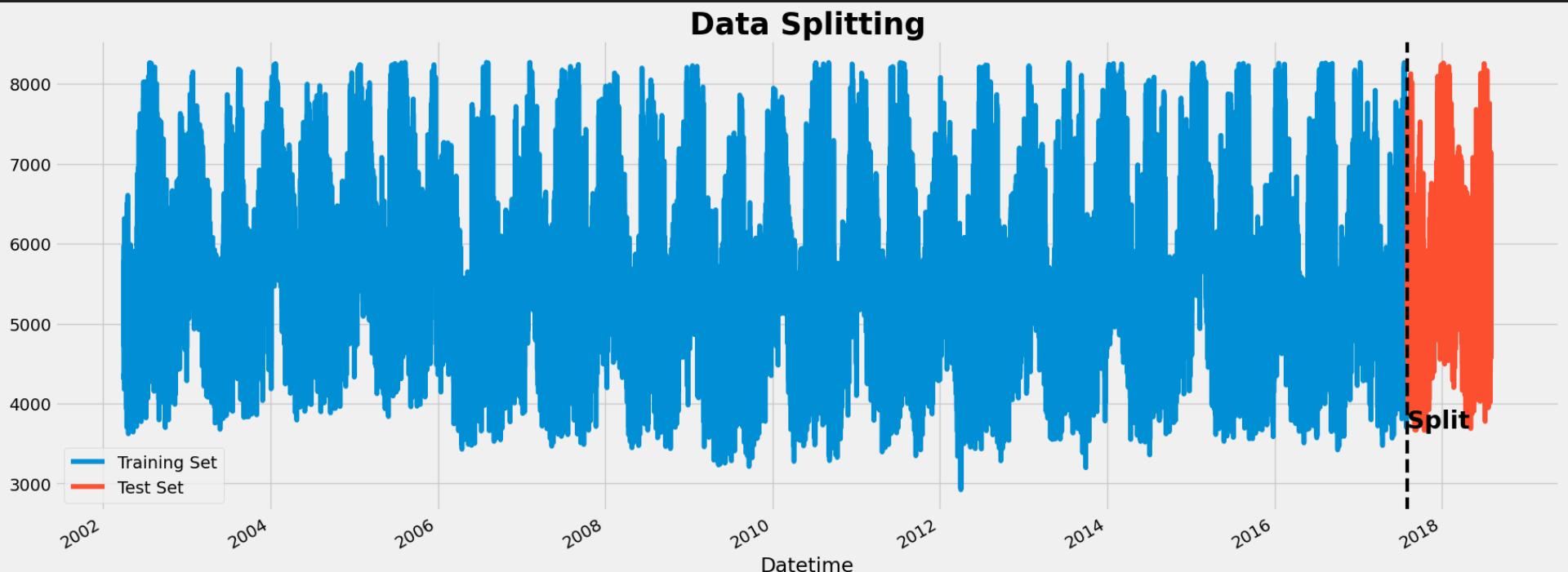
PHASE - 2

MODEL BUILDING

MODELS BUILT:

1. PROPHET MODEL
2. DEEP NEURAL NETWORK
3. XGBOOST
4. LINEAR REGRESSION
5. RANDOM FOREST

Splitting the data set

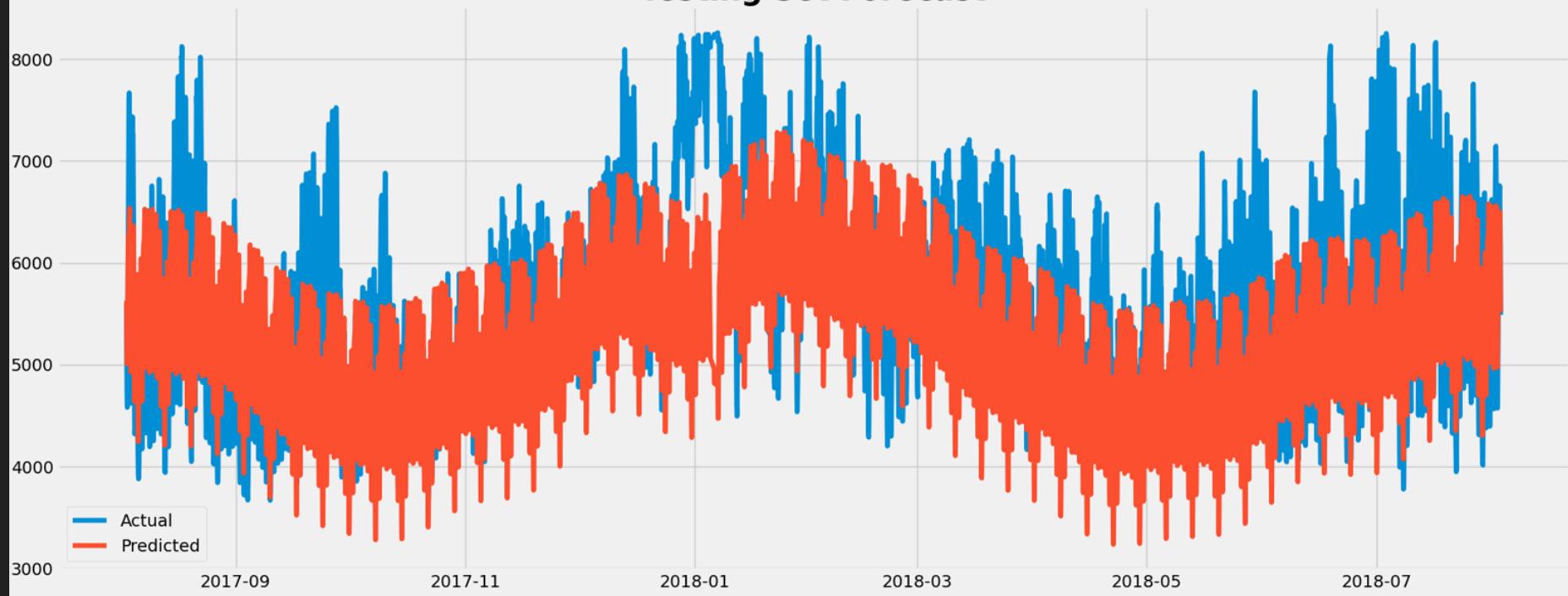


Train: 133852
Test: 8651

1. PROPHET MODEL

Prophet is an open-source library for univariate (one variable) time series forecasting developed by Facebook. It works best with time series that have strong seasonal effects and several seasons of historical data.

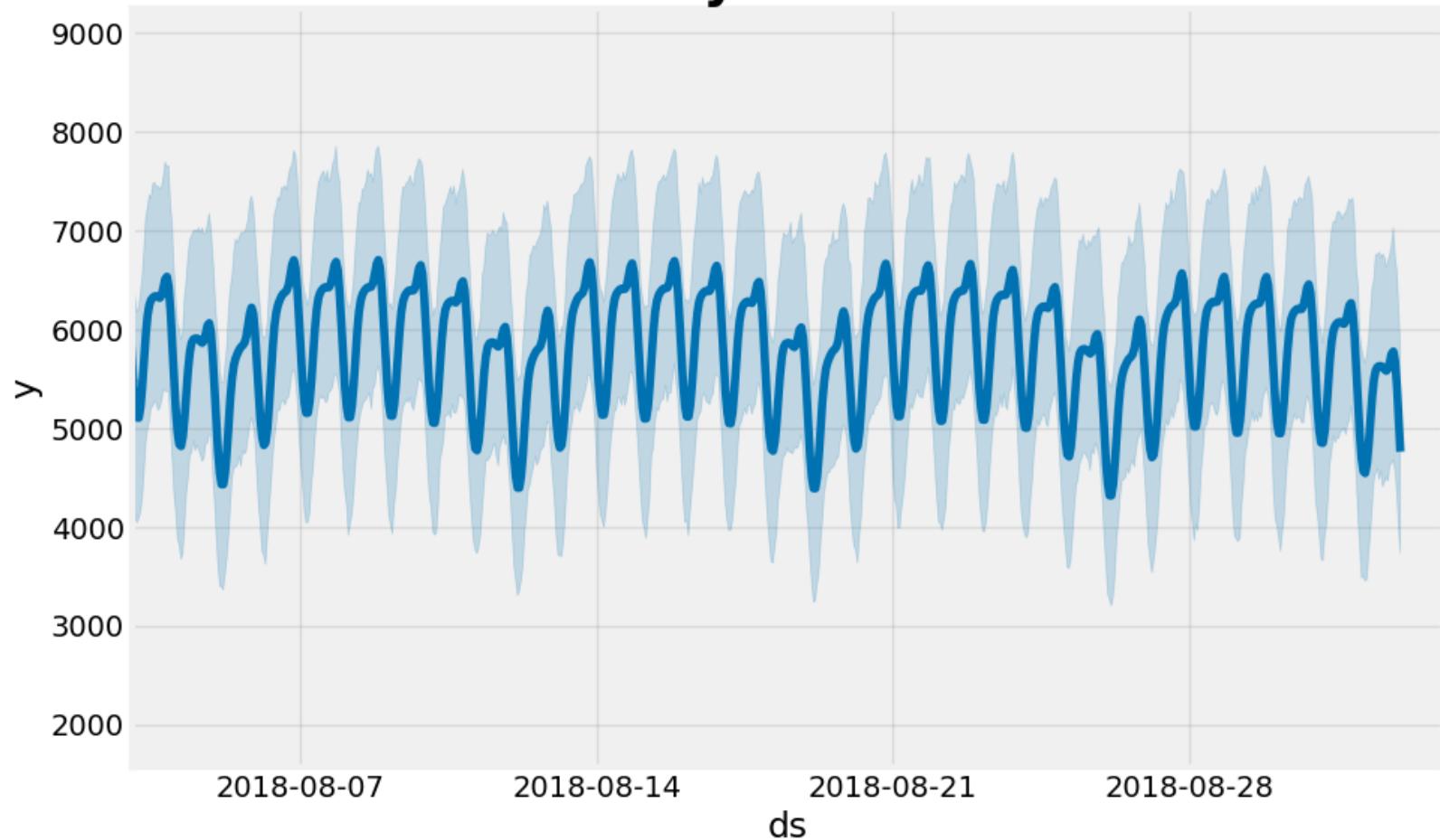
Testing Set Forecast

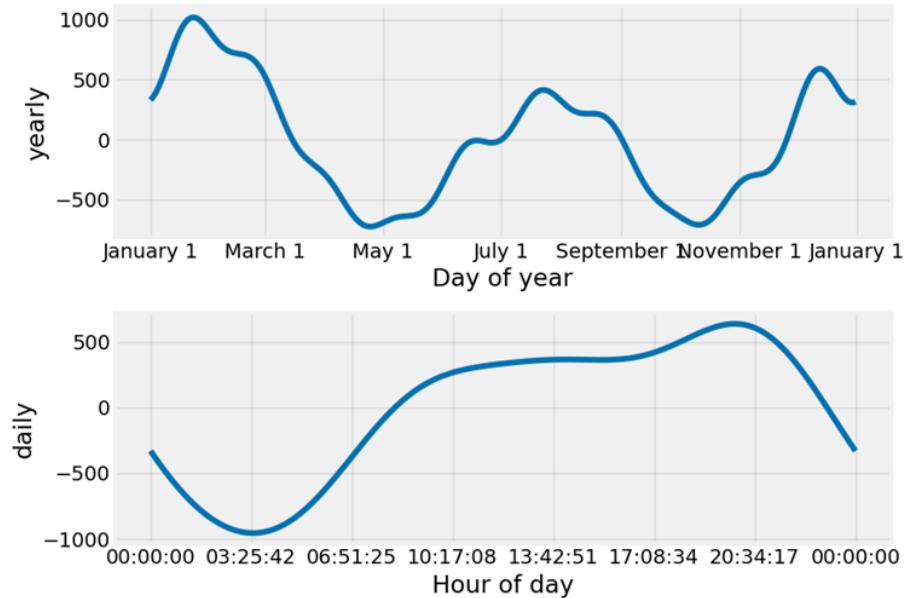
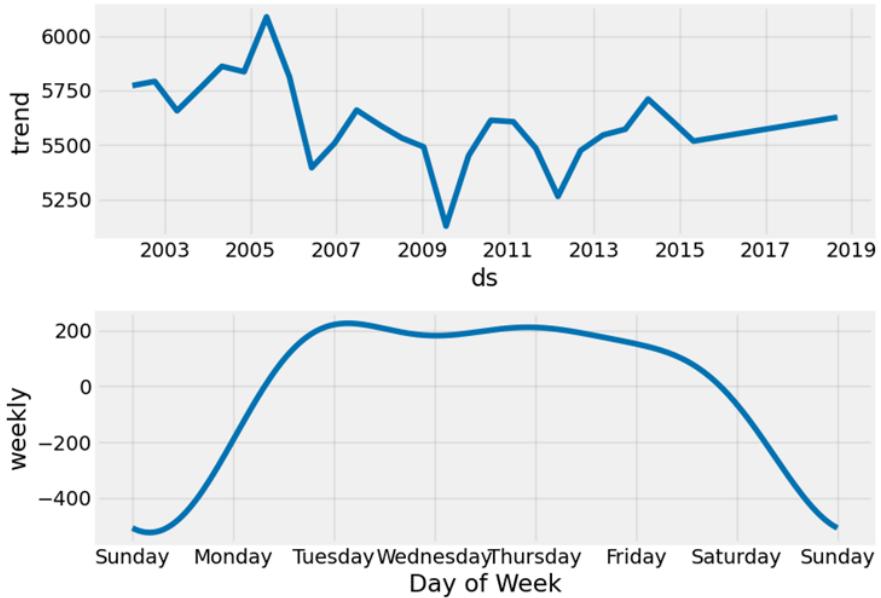


MAE: 545.605

RMSE Score on Test Set: 721.23

30 Days Forecast



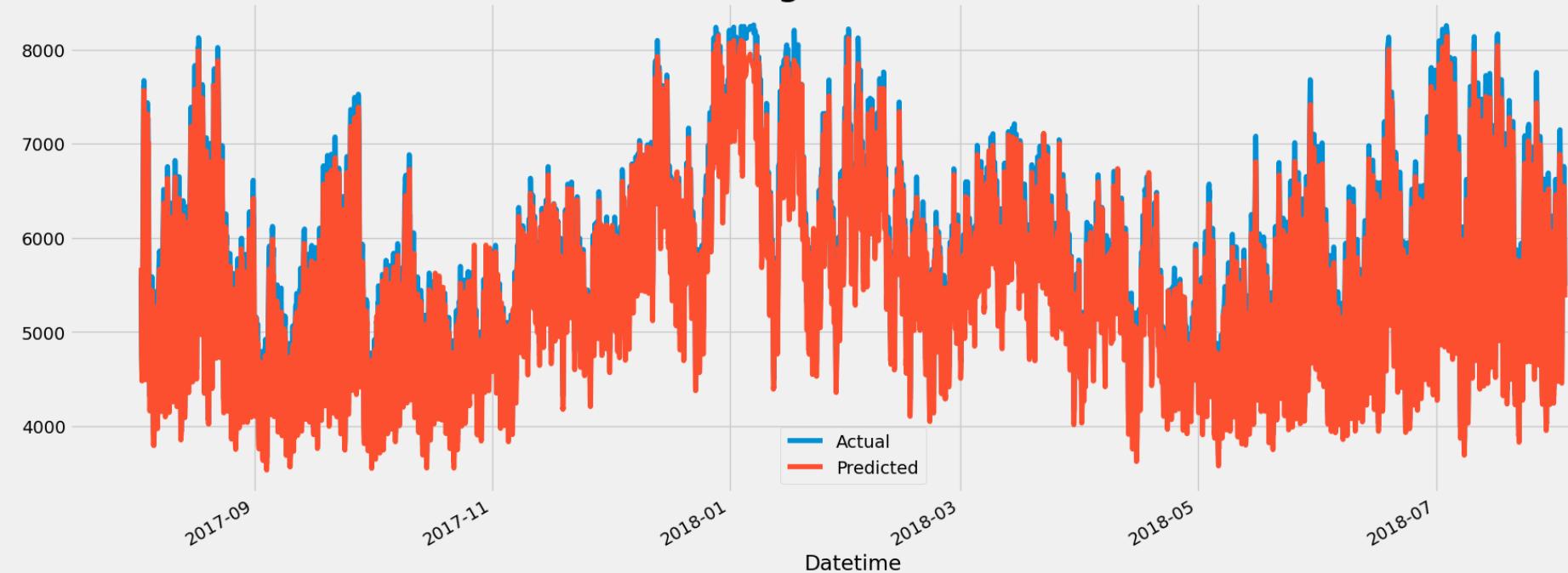


2. DEEP NEURAL NETWORK

The Deep Learning framework we are using is Tensorflow. Before feeding the data into Neural Network, we have to do some modification to the data so they can be accepted by the model.

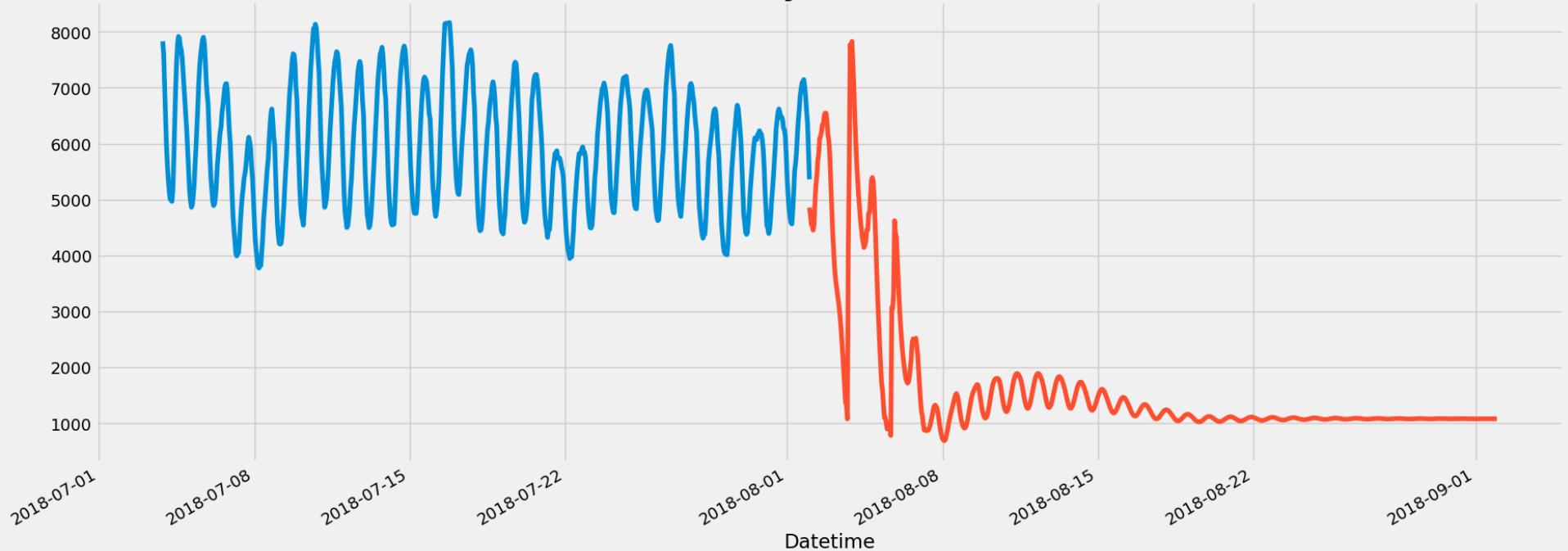
We are going to use windowing technique which basically group the data into feature and label. The label will be the next value. You can take a look at the next few cells to give an idea what we are going to do

Testing Set Forecast



MAE: 193.143

30 Days Forecast



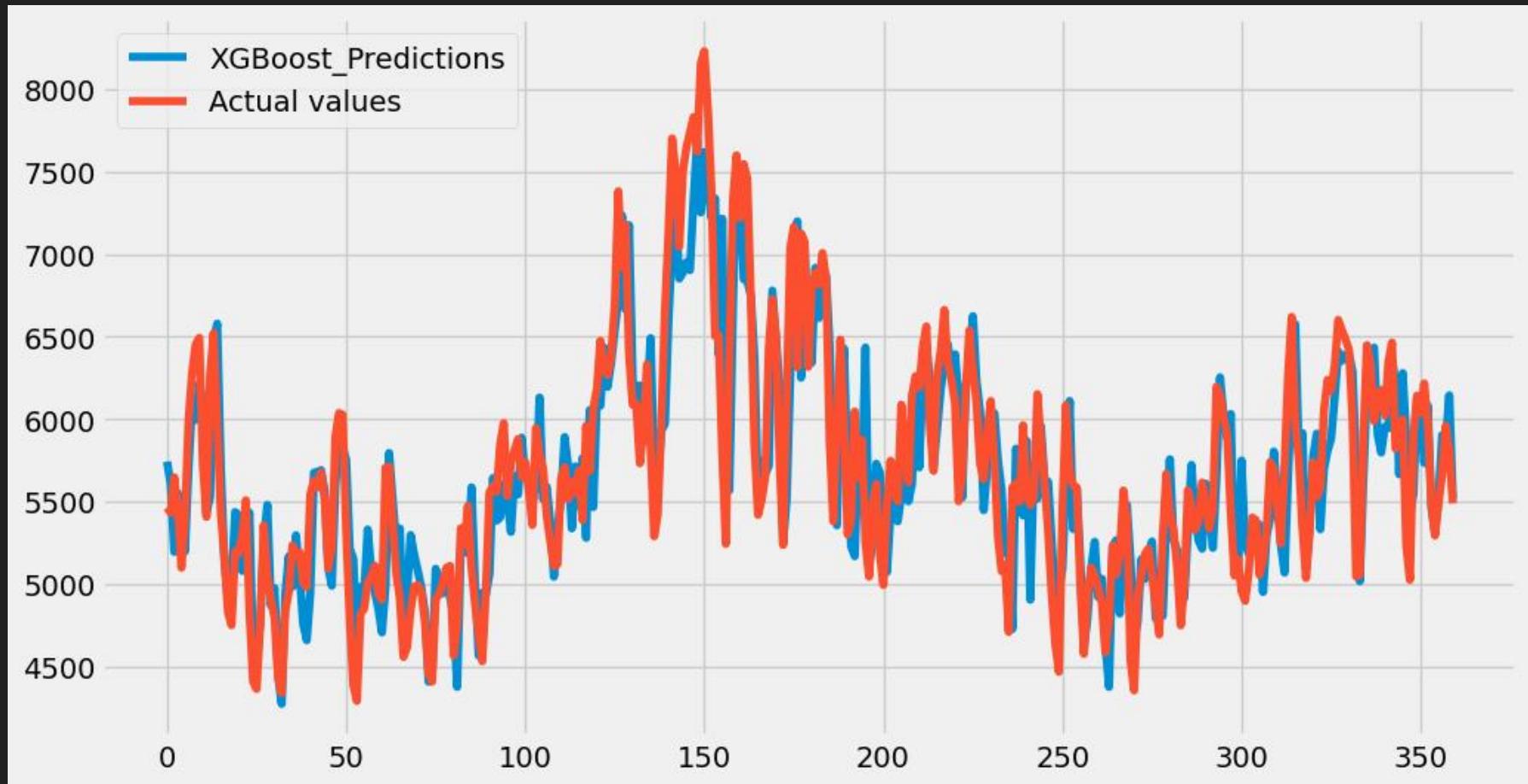
The Forecast of Deep Neural Network(DNN) model is not satisfactory.

3. XGBOOST MODEL

XGBoost is short for Extreme Gradient Boosting and is an efficient implementation of the stochastic gradient boosting machine learning algorithm. The stochastic gradient boosting algorithm, also called gradient boosting machines or tree boosting, is a powerful machine learning technique that performs well or even best on a wide range of challenging machine learning problems.

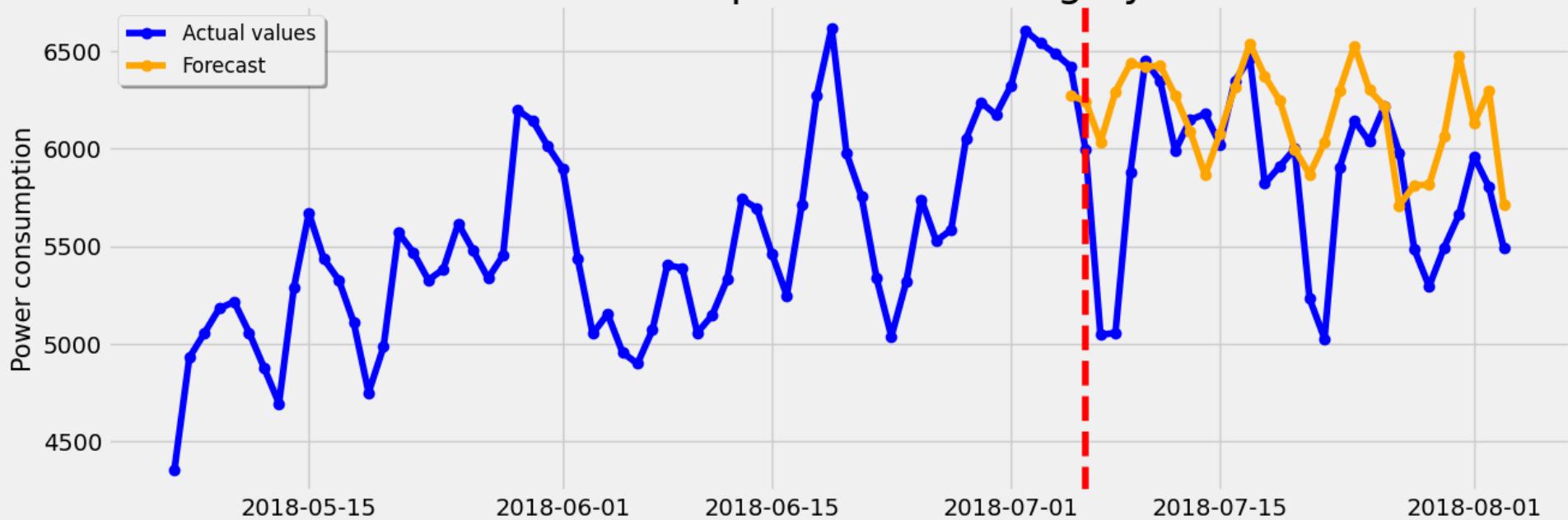
To be able to use XGBoost for time series forecasting, the data should be transformed into supervised learning before feeding it into the model.

Prediction has been done with test data



Forecasting

Power consumption forecasting by XGB



Predicted RMSE,MAE,MAPE

RMSE for XG Boost Model is: 353.39044761279615

MAE for XG Boos Model is: 270.524717518854

MAPE for XG Boos Model is: 4.708

Forecasted RMSE,MAE,MAPE

RMSE for XGB forecast check is: 485.2427674054524

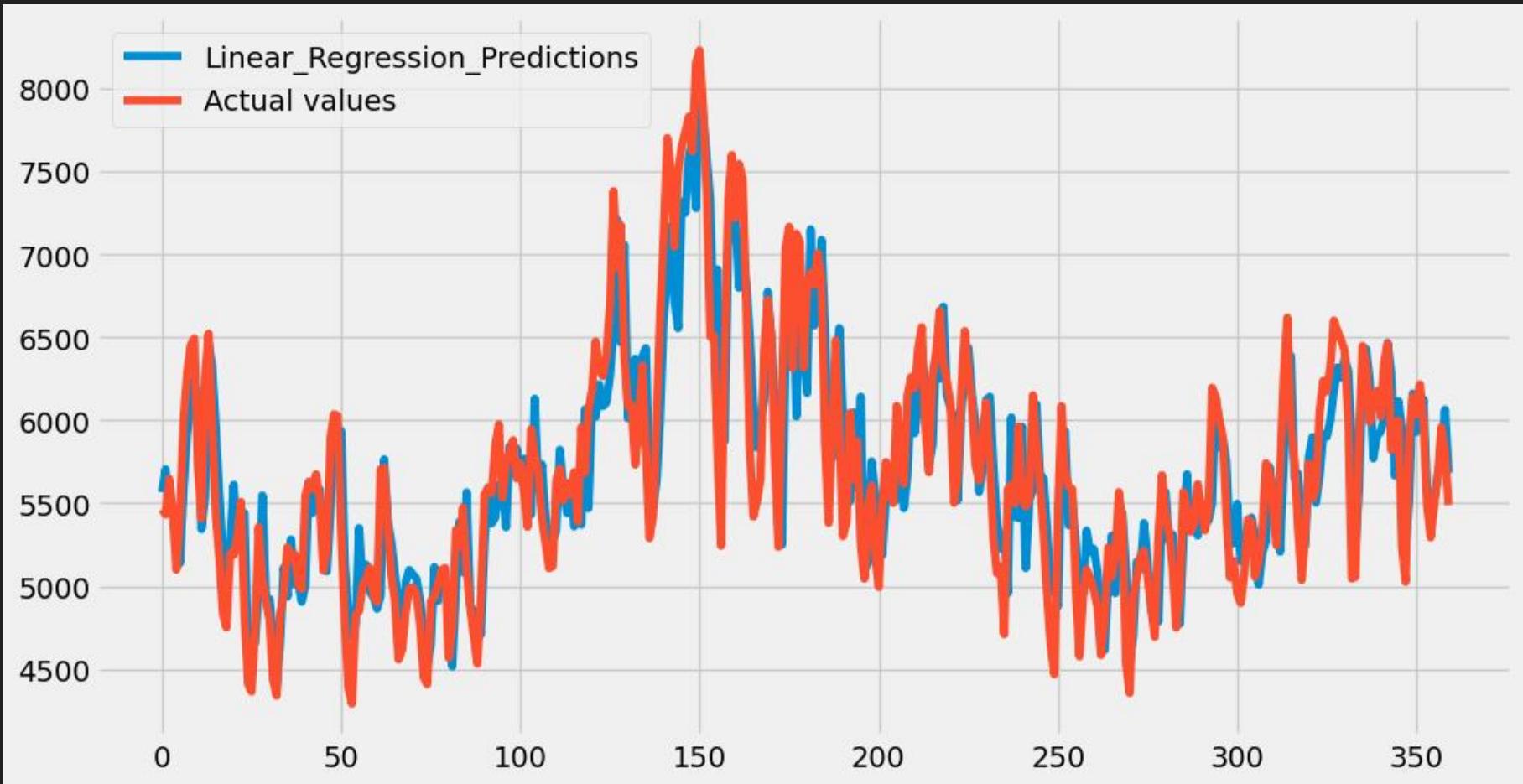
MAE for XGB forecast check is: 368.6199620225695

MAPE for XGB forecast check is: 6.668

4. LINEAR REGRESSION

Linear regression is a quiet and the simplest statistical regression method used for predictive analysis in machine learning. Linear regression shows the linear relationship between the independent(predictor) variable i.e. X-axis and the dependent(output) variable i.e. Y-axis, called linear regression. If there is a single input variable X(independent variable), such linear regression is called *simple linear regression*.

Prediction has been done with test data



Forecasting

Power consumption forecasting by LR



Predicted RMSE,MAE,MAPE

```
RMSE for Linear Regression Model is: 342.5745881638847  
MAE for Linear Regression Model is: 269.31156194432845  
MAPE for Linear Regression Model is: 4.703
```

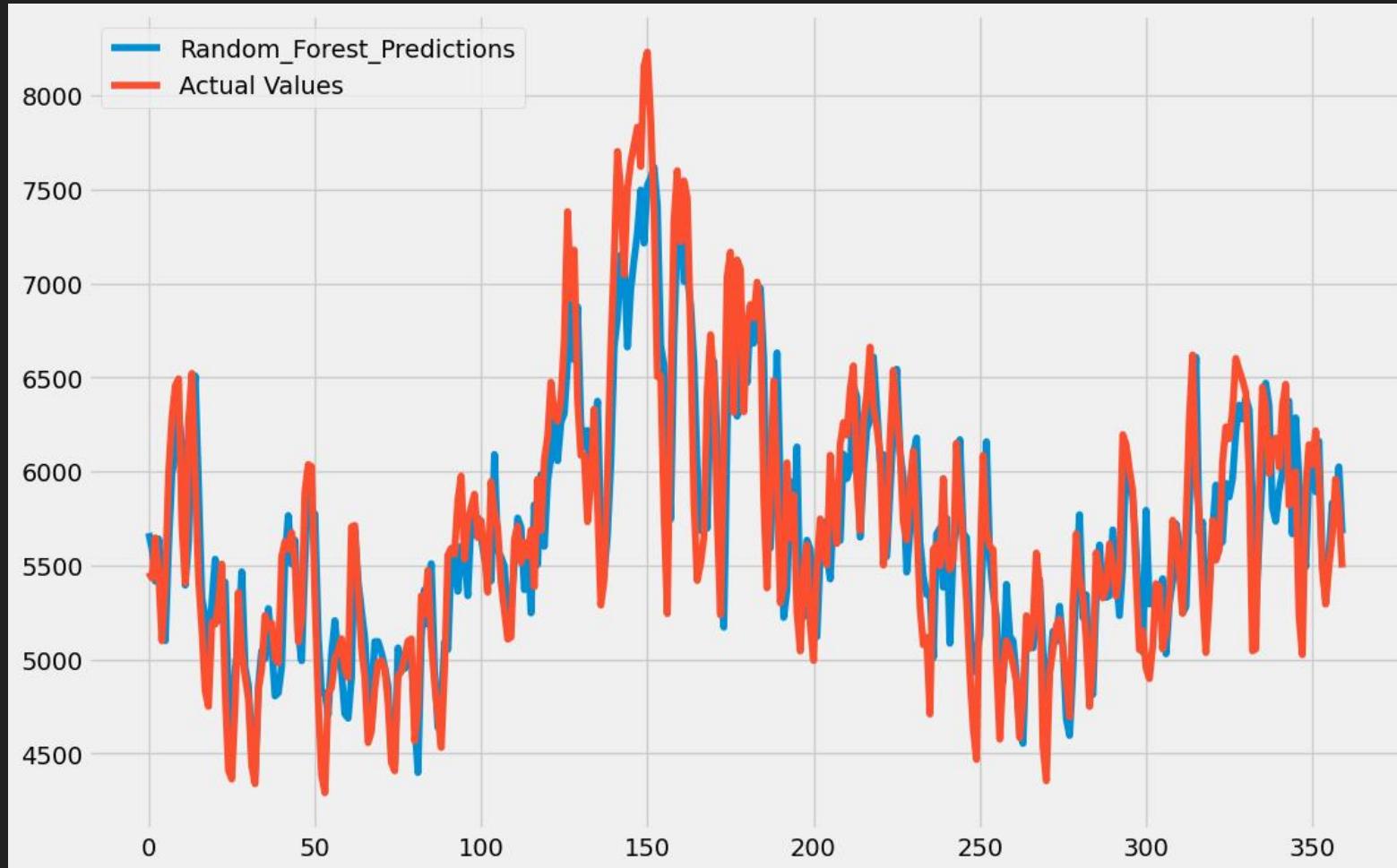
Forecasted RMSE,MAE,MAPE

```
RMSE for forecast check is: 748.1299433652413  
MAE for forecast check is: 604.5203961263184  
MAPE for forecast check is: 10.681
```

5. RANDOM FOREST

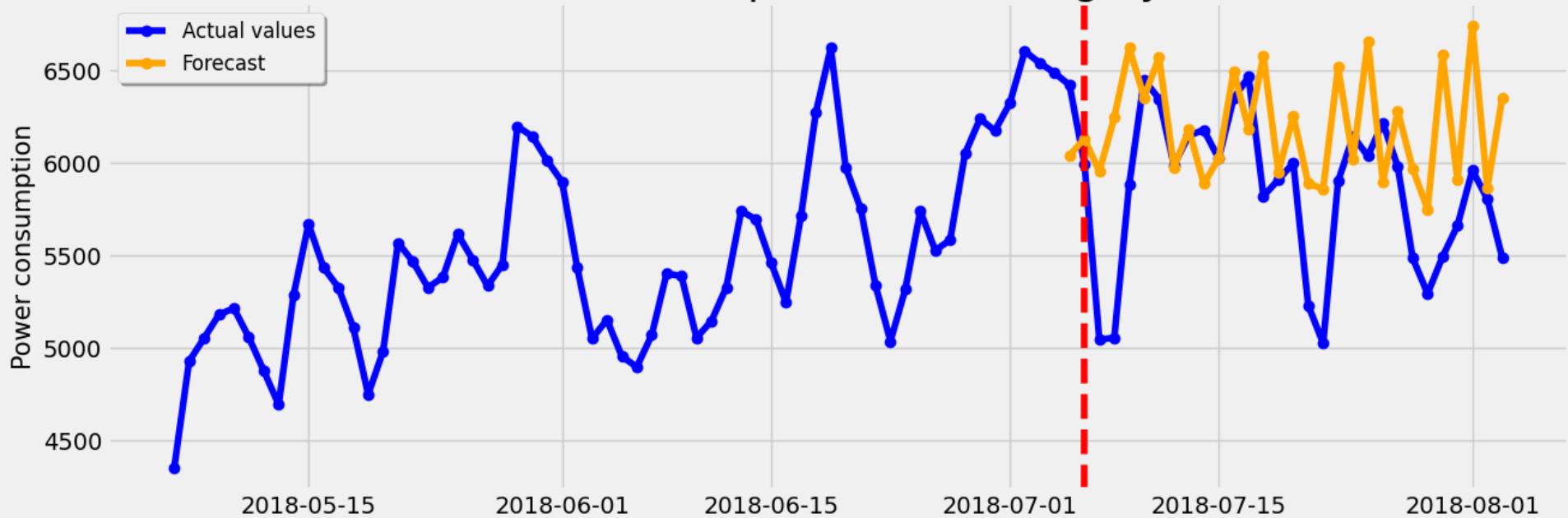
Random forest is a *Supervised Machine Learning Algorithm* that is *used widely in Classification and Regression problems*. It builds decision trees on different samples and takes their majority vote for classification and average in case of regression. One of the most important features of the Random Forest Algorithm is that it can handle the data set containing *continuous variables*, as in the case of regression, and *categorical variables*, as in the case of classification. It performs better for classification and regression tasks.

Prediction has been done with test data



Forecasting

Power consumption forecasting by RF



Predicted RMSE,MAE,MAPE

RMSE for Random Forest Model is: 346.4975114729343

MAE for Random Forest Model is: 273.41029119512467

MAPE for Random Forest Model is: 4.741

Forecasted RMSE,MAE,MAPE

RMSE for RF forecast check is: 546.8498682853016

MAE for RF forecast check is: 431.2064304612366

MAPE for RF forecast check is: 7.694

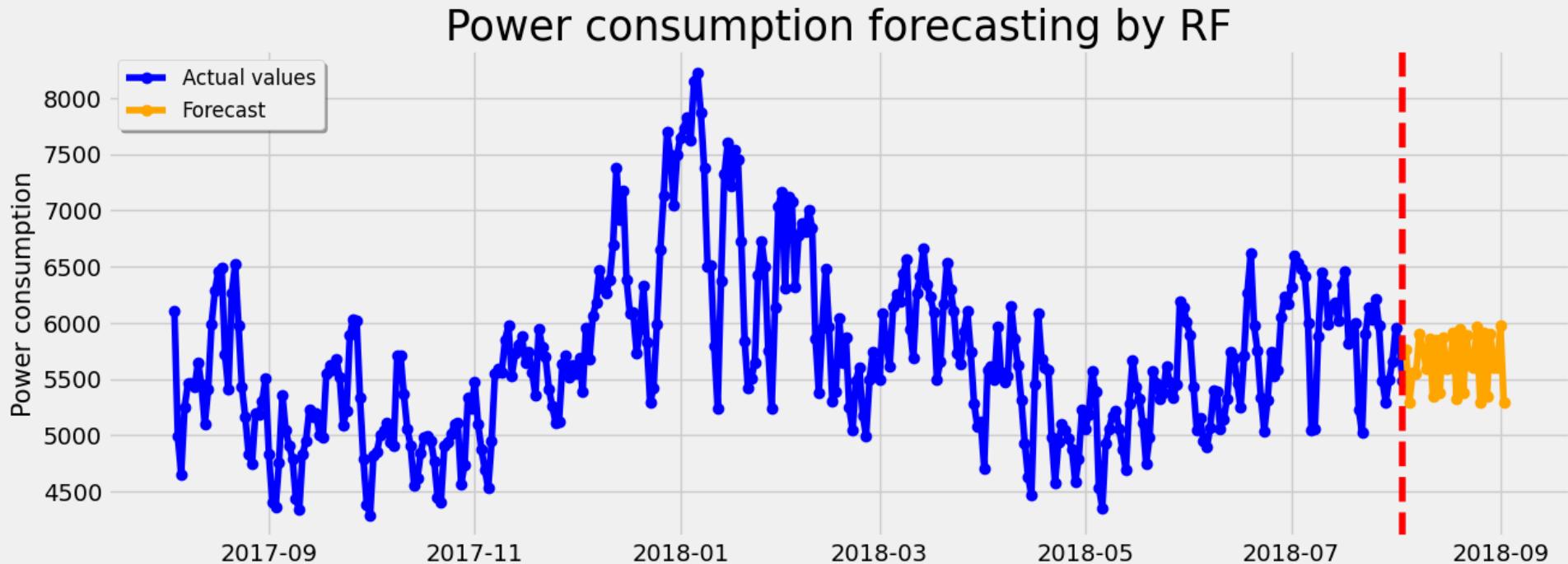
Findings

Based on our finding, we can see RANDOM FOREST model has less errors (RMSE, MAE, MAPE).

So forecasted power consumption for next 30 days from 2018-08-03 using RANDOM FOREST.

	future_energy
2018-08-04	5805.478750
2018-08-05	5250.396250
2018-08-06	5549.295942
2018-08-07	5510.736667
2018-08-08	5840.070000
2018-08-09	5759.415417
2018-08-10	5566.587138
2018-08-11	5990.728333
2018-08-12	5282.426250
2018-08-13	5849.890000
2018-08-14	5319.908333
2018-08-15	5782.970833
2018-08-16	5701.998333
2018-08-17	5539.374638
2018-08-18	6063.485417
2018-08-19	5293.225000
2018-08-20	5848.776667
2018-08-21	5237.801286
2018-08-22	5781.370000
2018-08-23	5712.072083
2018-08-24	5502.784167
2018-08-25	6005.816250
2018-08-26	5266.590833
2018-08-27	5932.142083
2018-08-28	5241.819583
2018-08-29	5781.027917
2018-08-30	5719.370833
2018-08-31	5523.631830
2018-09-01	6040.474583
2018-09-02	5221.987917

Forecasting for next 30 days

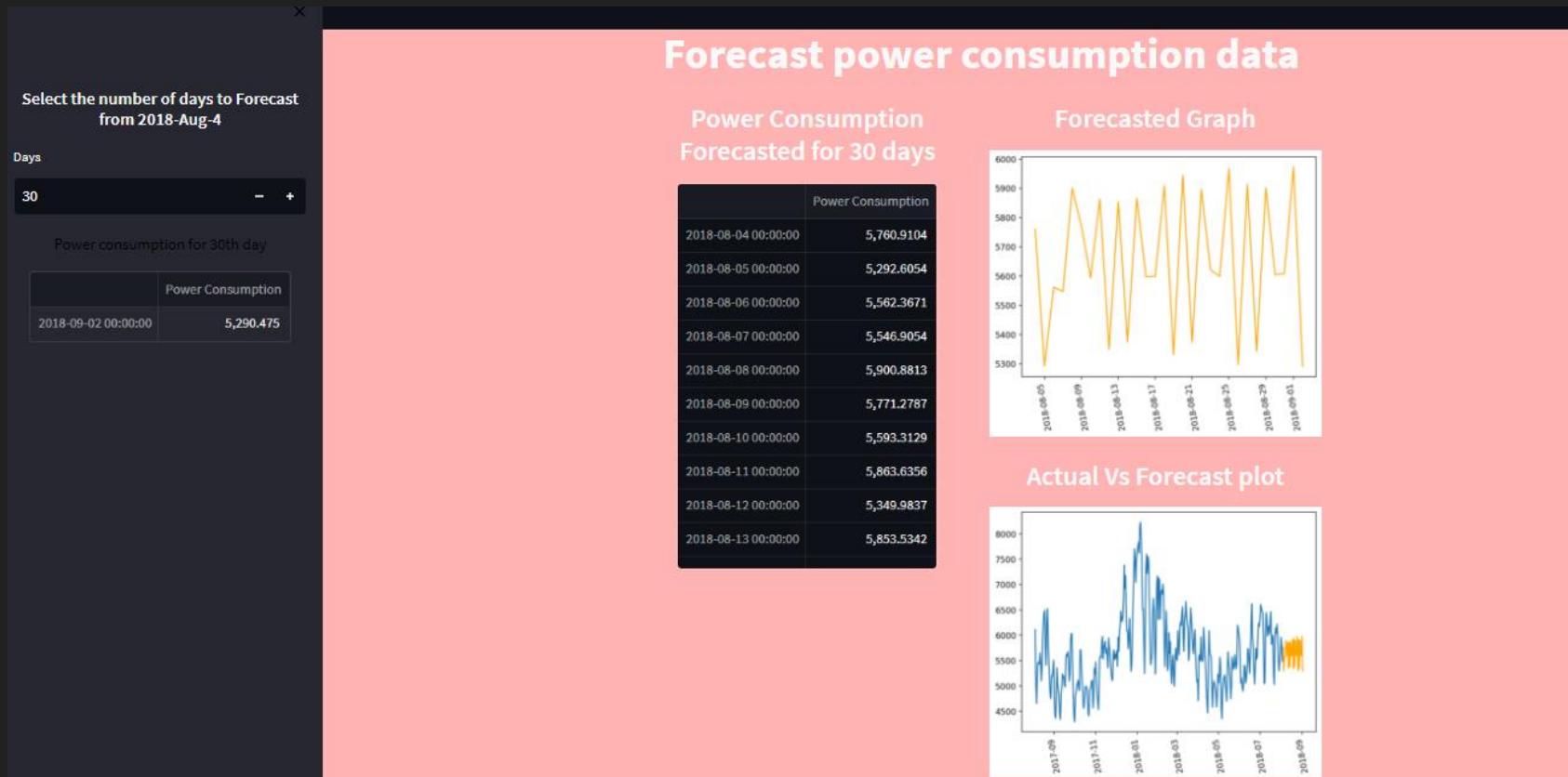


PHASE - 3

MODEL DEPLOYMENT

DEPLOYED MODEL :

The model is deployed on Streamlit framework and the Integrated Development Environment (IDE) used for its development was Spider.



Thank
You