

Introduction:

We will be taking a small forecasting problem and try to solve it till the end learning time series forecasting alongside.

What is Time Series analysis

Time series forecasting is a technique for the prediction of events through a sequence of time. The technique is used across many fields of study, from geology to behavior to economics. The techniques predict future events by [analyzing](#) the trends of the past, on the assumption that future trends will hold similar to historical trends.

Time series forecasting is performed in a variety of applications including:

- Weather forecasting
- Earthquake prediction
- Astronomy
- Statistics
- Mathematical finance
- Econometrics
- Pattern recognition
- [Signal](#) processing
- Control engineering

Time series forecasting is sometimes just the analysis of experts studying a field and offering their predictions. In many modern applications, however, time series forecasting uses computer technologies, including:

- [Machine learning](#)
- [Artificial neural networks](#)
- [Support vector machines](#)
- [Fuzzy logic](#)
- Gaussian processes
- Hidden [Markov models](#)

There are two main goals of time series analysis: (a) identifying the nature of the phenomenon represented by the sequence of observations, and (b) forecasting (predicting future values of the time series variable). Both of these goals require that the pattern of observed time series data is identified and more or less formally described. Once the pattern is established, we can interpret and integrate it with other data (i.e., use it in our theory of the investigated phenomenon, e.g., seasonal commodity prices). Regardless of the depth of our understanding and the validity of our interpretation (theory) of the phenomenon, we can extrapolate the identified pattern to predict future events.

Stages in Time Series Forecasting

Solving a time series problem is a little different as compared to a regular modeling task. A simple/basic journey of solving a time series problem can be demonstrated through the following processes. We will understand about tasks which one needs to perform in every stage. We will also look at the python implementation of each stage of our problem-solving journey.

Steps are –

1. Visualizing time series

In this step, we try to visualize the series. We try to identify all the underlying patterns related to the series like trend and seasonality. Do not worry about these terms right now, as we will discuss them during implementation. You can say that this is more a type of exploratory analysis of time series data.

2. Stationarising time series

A stationary time series is one whose statistical properties such as mean, variance, autocorrelation, etc. are all constant over time. Most statistical forecasting methods are based on the assumption that the time series can be rendered approximately stationary (i.e., “stationarised”) through the use of mathematical transformations. A stationarised series is relatively easy to predict: you simply predict that its statistical properties will be the same in the future as they have been in the past! Another reason for trying to stationarise a time series is to be able to obtain meaningful sample statistics such as means, variances, and correlations with other variables. Such statistics are useful as descriptors of future behavior only if the series is stationary. For example, if the series is consistently increasing over time, the sample mean and variance will grow with the size of the sample, and they will always underestimate the mean and variance in future periods. And if the mean and variance of a series are not well-defined, then neither are its correlations with other variables.

3. Finding the best parameters for our model

We need to find optimal parameters for forecasting models once we have a stationary series. These parameters come from the ACF and PACF plots. Hence, this stage is more about plotting above two graphs and extracting optimal model parameters based on them. Do not worry, we will cover on how to determine these parameters during the implementation part below!

4. Fitting model

Once we have our optimal model parameters, we can fit an ARIMA model to learn the pattern of the series. Always remember that time series algorithms work on stationary data only hence making a series stationary is an important aspect

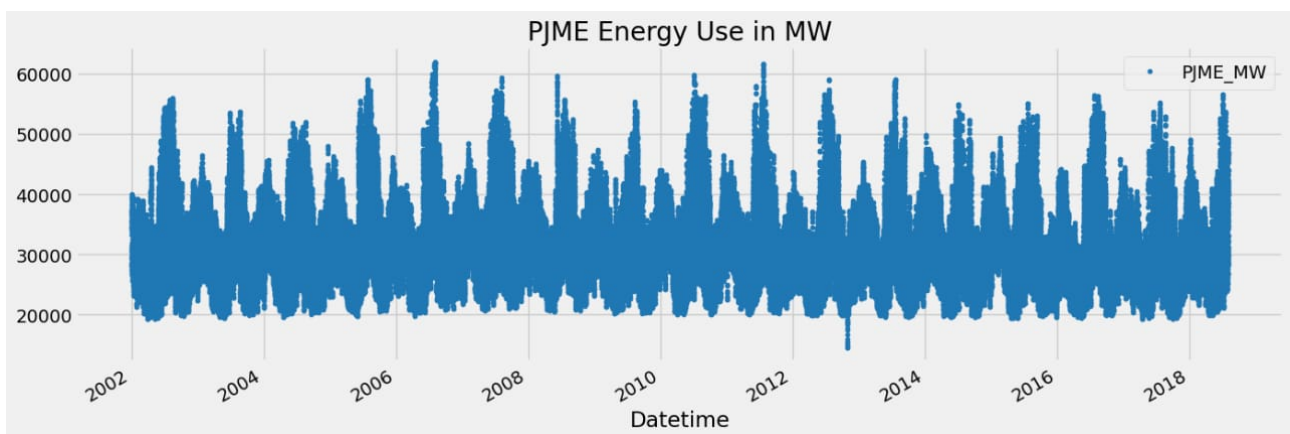
5. Predictions

After fitting our model, we will be predicting the future in this stage. Since we are now familiar with a basic flow of solving a time series problem, let us get to the implementation.

Problem Statement

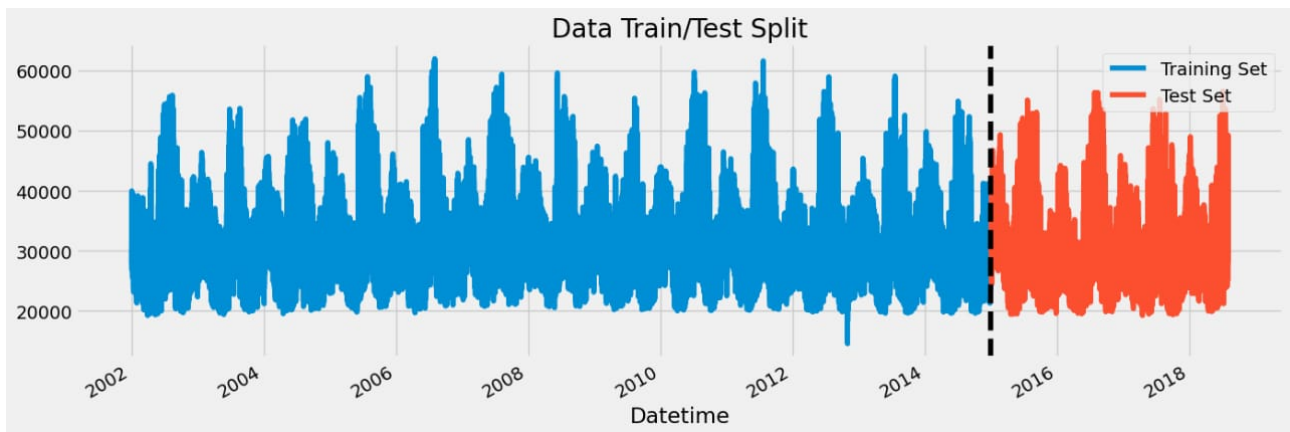
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import xgboost as xgb
from sklearn.metrics import mean_squared_error
color_pal = sns.color_palette()
plt.style.use('fivethirtyeight')
df = pd.read_csv('../input/hourly-energy-consumption/PJME_hourly.csv')
df = df.set_index('Datetime')
df.index = pd.to_datetime(df.index)
df.plot(style='.',
        figsize=(15, 5),
        color=color_pal[0],
        title='PJME Energy Use in MW')
plt.show()
```

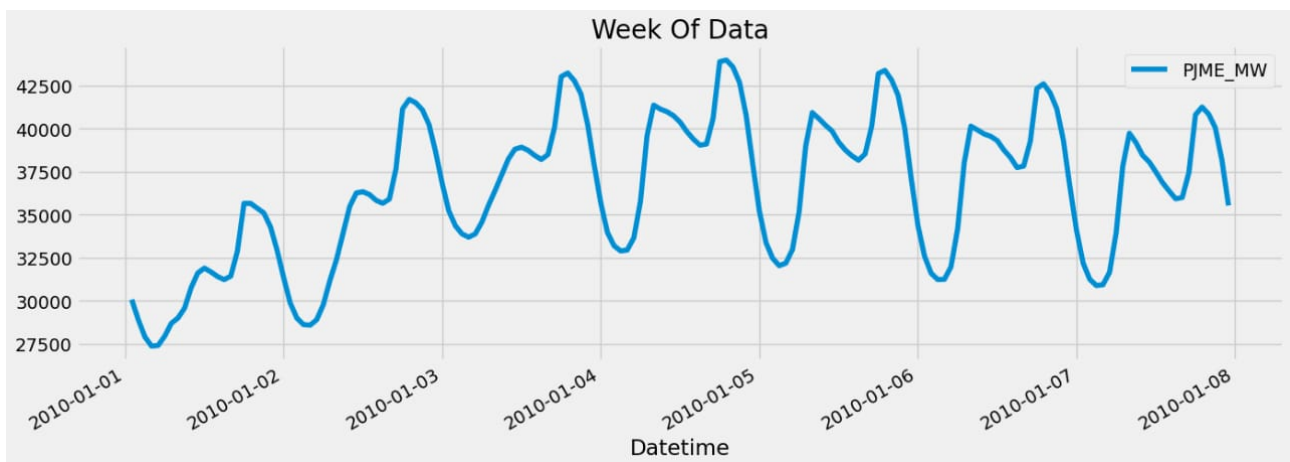


```
train = df.loc[df.index < '01-01-2015']
test = df.loc[df.index >= '01-01-2015']
```

```
fig, ax = plt.subplots(figsize=(15, 5))
train.plot(ax=ax, label='Training Set', title='Data Train/Test Split')
test.plot(ax=ax, label='Test Set')
ax.axvline('01-01-2015', color='black', ls='--')
ax.legend(['Training Set', 'Test Set'])
plt.show()
```



```
df.loc[(df.index > '01-01-2010') & (df.index < '01-08-2010')] \
.plot(figsize=(15, 5), title='Week Of Data')
plt.show()
```

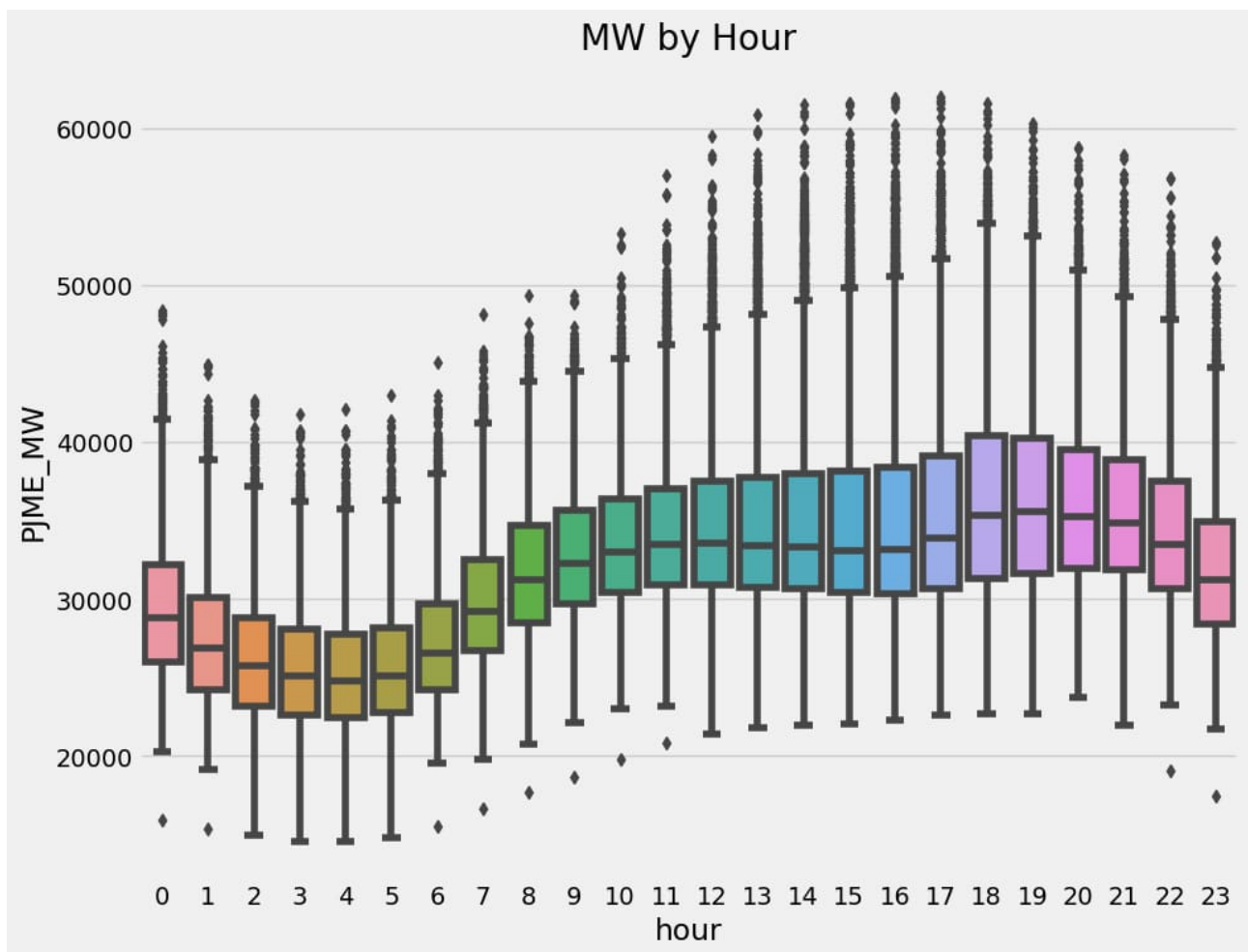


```
def create_features(df):
    """
    Create time series features based on time series index.
    """
```

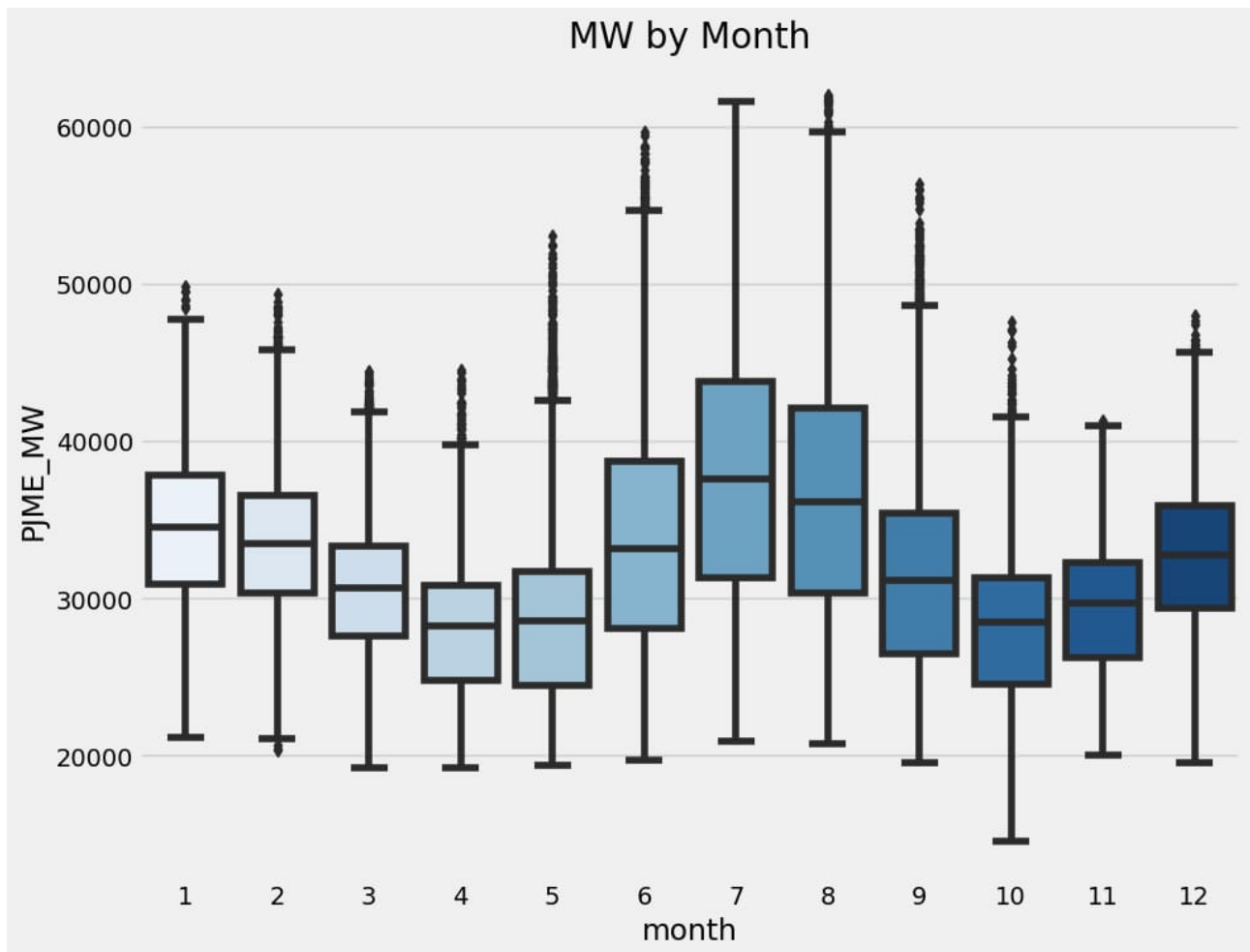
```
    df = df.copy()
    df['hour'] = df.index.hour
    df['dayofweek'] = df.index.dayofweek
    df['quarter'] = df.index.quarter
    df['month'] = df.index.month
    df['year'] = df.index.year
    df['dayofyear'] = df.index.dayofyear
    df['weekofmonth'] = df.index.day
    df['weekofyear'] = df.index.isocalendar().week
    return df
```

```
df = create_features(df)
```

```
fig, ax = plt.subplots(figsize=(10, 8))
sns.boxplot(data=df, x='hour', y='PJME_MW')
ax.set_title('MW by Hour')
plt.show()
```



```
fig, ax = plt.subplots(figsize=(10, 8))
sns.boxplot(data=df, x='month', y='PJME_MW', palette='Blues')
ax.set_title('MW by Month')
plt.show()
```



```
train = create_features(train)
```

```
test = create_features(test)
```

```
FEATURES = ['dayofyear', 'hour', 'dayofweek', 'quarter', 'month', 'year']
```

```
TARGET = 'PJME_MW'
```

```
x_train = train[FEATURES]
```

```
y_train = train[TARGET]
```

```
x_test = test[FEATURES]
```

```
y_test = test[TARGET]
```

```
reg = xgb.XGBRegressor(base_score=0.5, booster='gbtree', n_estimators=1000,
                        early_stopping_rounds=50,
                        objective='reg:linear',
                        max_depth=3,
                        learning_rate=0.01)
```

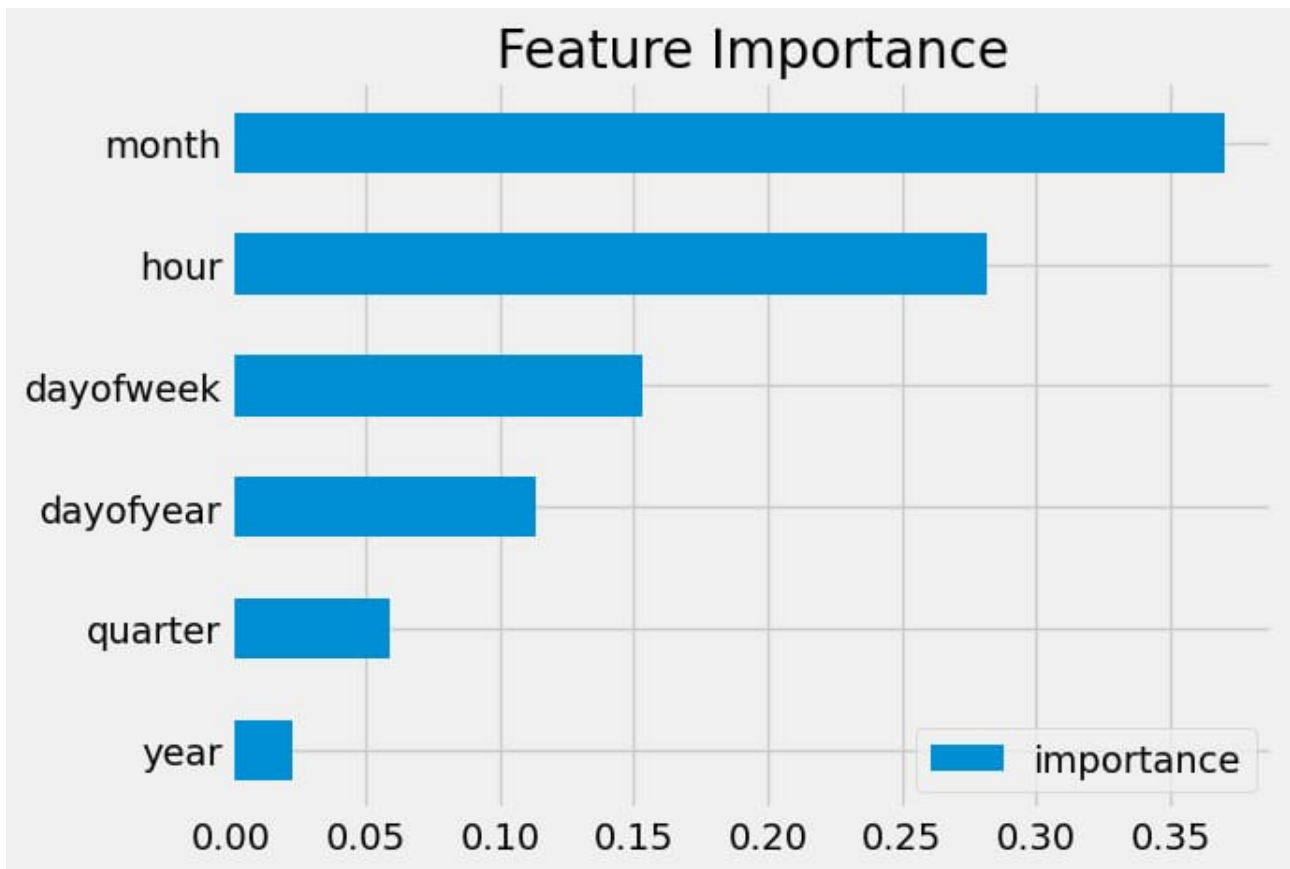
```
reg.fit(x_train, y_train,
        eval_set= [(x_train, y_train), (x_test, y_test)],
        verbose=100)
```

```
fi = pd.DataFrame(data=reg.feature_importances_,
                  index=reg.feature_names_in_,
```

```

columns=['importance'])
fi.sort_values('importance').plot(kind='barh', title='Feature Importance')
plt.show()

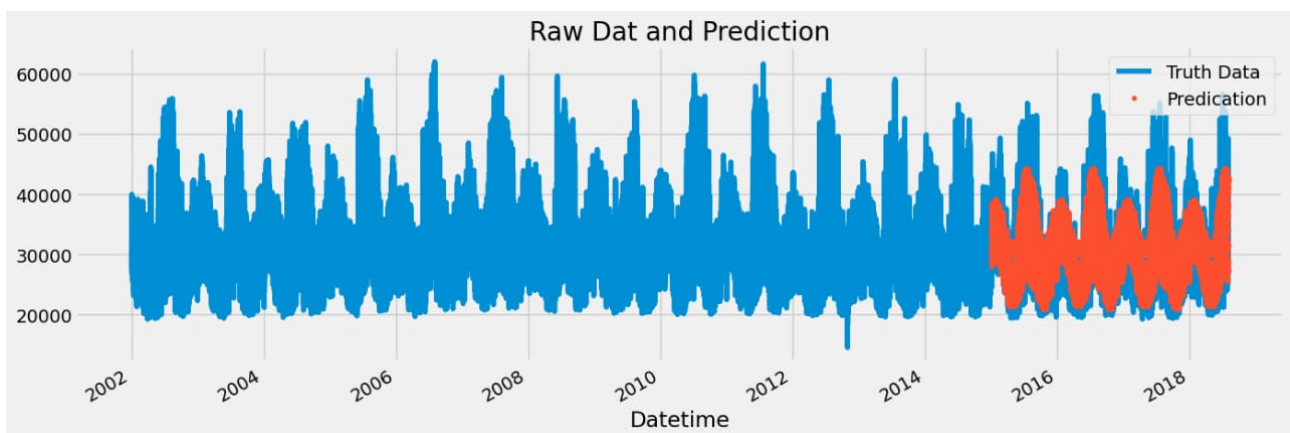
```



```

test['prediction'] = reg.predict(x_test)
df = df.merge(test[['prediction']], how= 'left', left_index=True, right_index=True)
ax = df[['PJME_MW']].plot(figsize=(15, 5))
df['prediction'].plot(ax=ax, style='.')
plt.legend(['Truth Data', 'Predication'])
ax.set_title('Raw Dat and Prediction')
plt.show()

```

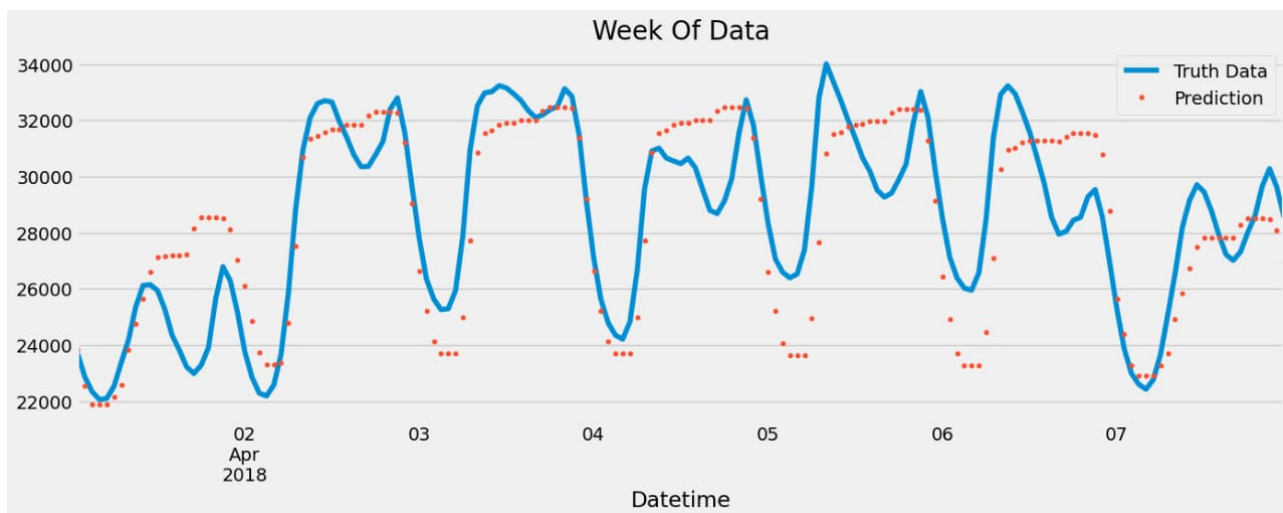


```

ax = df.loc[(df.index > '04-01-2018') & (df.index < '04-08-2018')][['PJME_MW']] \
    .plot(figsize=(15, 5), title='Week Of Data')
df.loc[(df.index > '04-01-2018') & (df.index < '04-08-2018')][['prediction']] \
    .plot(style='.')

```

```
plt.legend(['Truth Data','Prediction'])
plt.show()
```



```
score = np.sqrt(mean_squared_error(test['PJME_MW'], test['prediction']))
```

```
print(f'RMSE Score on Test set: {score:0.2f}')
```

RMSE Score on Test set: 3721.75

Calculate Error Look at the worst and best predicted days

```
test['error'] = np.abs(test[TARGET] - test['prediction'])
```

```
test['date'] = test.index.date
```

```
test.groupby(['date'])['error'].mean().sort_values(ascending=False).head(10)
```

```
date
```

Output:

2016-08-13 12839.597087

2016-08-14 12780.209961

2016-09-10 11356.302979

2015-02-20 10965.982259

2016-09-09 10864.954834

2018-01-06 10506.845622

2016-08-12 10124.051595

2015-02-21 9881.803711

2015-02-16 9781.552246

2018-01-07 9739.144206

Name: error, dtype: float64