



# Python for TIME SERIES FORECASTING

Using Machine Learning

ARIMA

SARIMAX

ARIMA

LSTM &  
Prophet

DATA PREPROCESSING



## EVALUATION METRICS

- MAE
- RMSE



Predicting the Future



## ARIMA Times Forecasting Workflow



# Time Series Analysis and Forecasting with ARIMA

Load the data

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Load the data
data = pd.read_csv('TimeSeries_TotalSolarGen_and_Load_IT_2016.csv')
data.head()
```

```
Out[1]:    utc_timestamp  IT_load_new  IT_solar_generation
0  2016-01-01T00:00:00Z      21665.0                  1
1  2016-01-01T01:00:00Z      20260.0                  0
2  2016-01-01T02:00:00Z      19056.0                  0
3  2016-01-01T03:00:00Z      18407.0                  0
4  2016-01-01T04:00:00Z      18425.0                  0
```

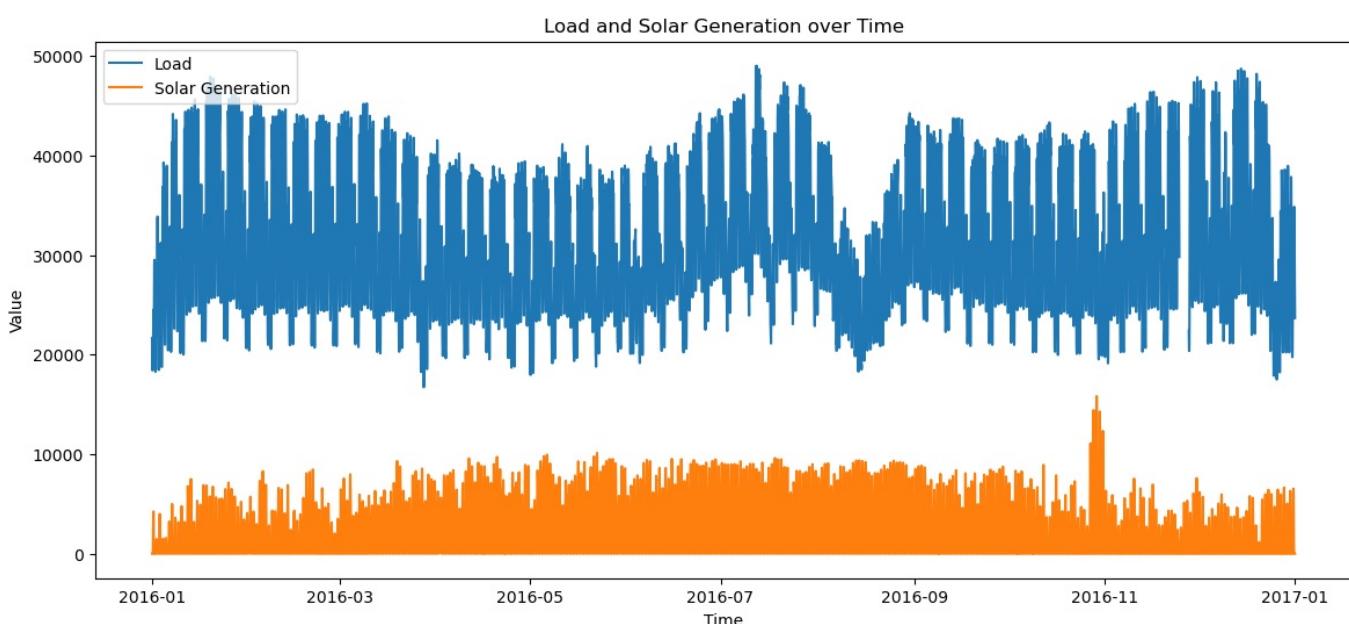
## Visualize the data

```
In [2]: import matplotlib.pyplot as plt

# Convert utc_timestamp to datetime
data['utc_timestamp'] = pd.to_datetime(data['utc_timestamp'])

# Plot the data
plt.figure(figsize=(14,6))

plt.plot(data['utc_timestamp'], data['IT_load_new'], label='Load')
plt.plot(data['utc_timestamp'], data['IT_solar_generation'], label='Solar Generation')
plt.xlabel('Time')
plt.ylabel('Value')
plt.legend()
plt.title('Load and Solar Generation over Time')
plt.show()
```



```
In [3]: import matplotlib.pyplot as plt

# Convert utc_timestamp to datetime
data['utc_timestamp'] = pd.to_datetime(data['utc_timestamp'])

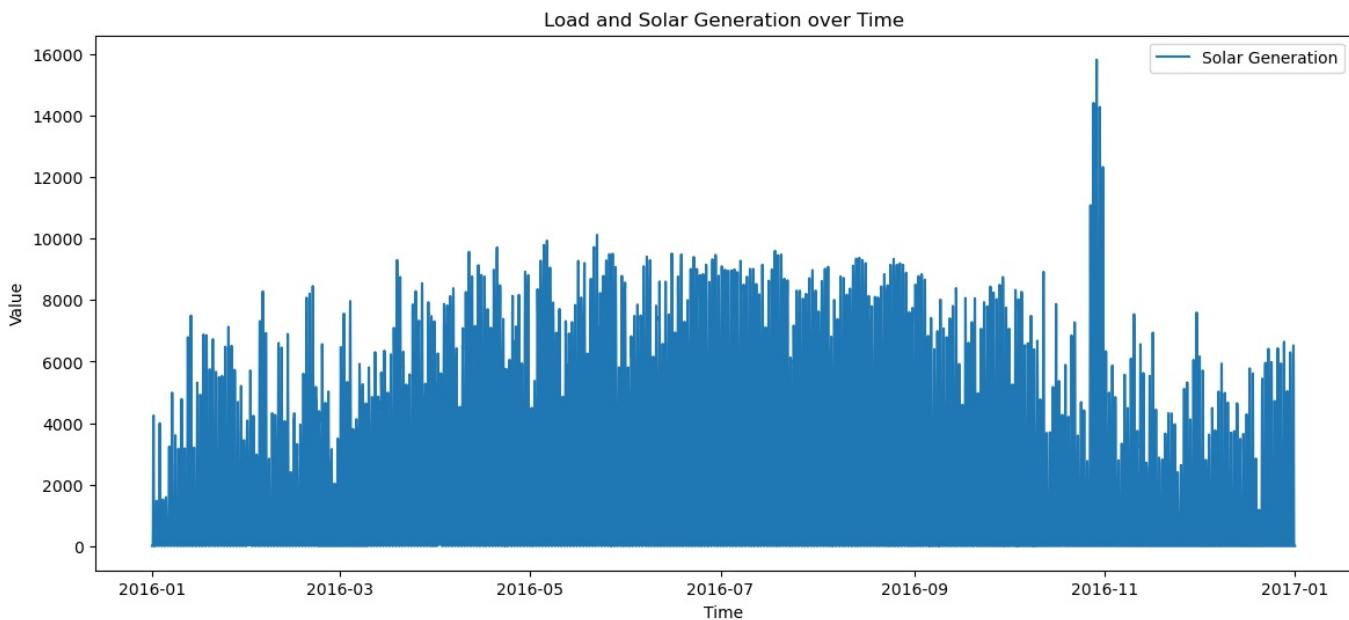
# Plot the data
plt.figure(figsize=(14,6))

# plt.plot(data['utc_timestamp'], data['IT_load_new'], label='Load')
plt.plot(data['utc_timestamp'], data['IT_solar_generation'], label='Solar Generation')
```

```

plt.xlabel('Time')
plt.ylabel('Value')
plt.legend()
plt.title('Load and Solar Generation over Time')
plt.show()

```



The plot shows both the load and solar generation over time for the year 2016. The load appears to have a cyclical pattern with peaks and valleys, possibly corresponding to daily patterns of electricity use.

The solar generation also shows a clear pattern, with generation during the day and no generation at night (as expected). The amount of solar generation also appears to fluctuate throughout the year, likely due to seasonal changes in sunlight.

Before moving on to time series analysis and forecasting with ARIMA, it's important to check for stationarity in your time series data. Stationarity is a property of time series data that implies the mean, variance, and autocorrelation structure do not change over time. Many time series models, including ARIMA, require the data to be stationary.

Let's perform an Augmented Dickey-Fuller test to check the stationarity of the time series. The null hypothesis of the ADF test is that the time series is non-stationary. So, if the p-value of the test is less than the significance level (0.05) then you reject the null hypothesis and infer that the time series is indeed stationary.

## Handle missing value

```
In [4]: # Check for missing values
data.isnull().sum()
```

```
Out[4]: utc_timestamp      0
IT_load_new        72
IT_solar_generation      0
dtype: int64
```

```
In [5]: # Fill missing values using forward fill
data['IT_load_new'].fillna(method='ffill', inplace=True)

# Check for missing values again
print("Missing values after filling:")
print(data.isnull().sum())
```

```
Missing values after filling:
utc_timestamp      0
IT_load_new        0
IT_solar_generation      0
dtype: int64
```

```
C:\Users\asus1\AppData\Local\Temp\ipykernel_33888\3163531581.py:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.  
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.
```

```
For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.
```

```
data['IT_load_new'].fillna(method='ffill', inplace=True)  
C:\Users\asus1\AppData\Local\Temp\ipykernel_33888\3163531581.py:2: FutureWarning: Series.fillna with 'method' is deprecated and will raise in a future version. Use obj.ffmpeg() or obj.bfill() instead.  
data['IT_load_new'].fillna(method='ffill', inplace=True)
```

## Check for stationarity

```
In [6]: from statsmodels.tsa.stattools import adfuller  
  
# Function to perform Augmented Dickey-Fuller test  
def adf_test(timeseries):  
    print ('Results of Dickey-Fuller Test:')    dfoutput = adfuller(timeseries, autolag='AIC')  
    dfoutput = pd.Series(dfoutput[0:4], index=['Test Statistic','p-value','#Lags Used','Number of Observations Used'])  
    for key,value in dfoutput[4].items():  
        dfoutput['Critical Value (%s)'%key] = value  
    print (dfoutput)
```

```
In [7]: # Perform Augmented Dickey-Fuller test again  
print("\nADF test for 'IT_load_new' after filling missing values:")  
adf_test(data['IT_load_new'])  
  
print("\nADF test for 'IT_solar_generation':")  
adf_test(data['IT_solar_generation'])
```

```
ADF test for 'IT_load_new' after filling missing values:  
Results of Dickey-Fuller Test:  
Test Statistic           -1.197390e+01  
p-value                  3.841445e-22  
#Lags Used                3.700000e+01  
Number of Observations Used 8.746000e+03  
Critical Value (1%)      -3.431098e+00  
Critical Value (5%)       -2.861871e+00  
Critical Value (10%)      -2.566946e+00  
dtype: float64
```

```
ADF test for 'IT_solar_generation':  
Results of Dickey-Fuller Test:  
Test Statistic           -5.741335e+00  
p-value                  6.265438e-07  
#Lags Used                3.600000e+01  
Number of Observations Used 8.747000e+03  
Critical Value (1%)      -3.431098e+00  
Critical Value (5%)       -2.861870e+00  
Critical Value (10%)      -2.566946e+00  
dtype: float64
```

```
In [8]: # Perform test for 'IT_load_new'  
print("ADF test for 'IT_load_new':")  
adf_test(data['IT_load_new'])  
  
# Perform test for 'IT_solar_generation'  
print("\nADF test for 'IT_solar_generation':")  
adf_test(data['IT_solar_generation'])
```

```

ADF test for 'IT_load_new':
Results of Dickey-Fuller Test:
Test Statistic           -1.197390e+01
p-value                  3.841445e-22
#Lags Used              3.700000e+01
Number of Observations Used 8.746000e+03
Critical Value (1%)      -3.431098e+00
Critical Value (5%)       -2.861871e+00
Critical Value (10%)      -2.566946e+00
dtype: float64

```

```

ADF test for 'IT_solar_generation':
Results of Dickey-Fuller Test:
Test Statistic           -5.741335e+00
p-value                  6.265438e-07
#Lags Used              3.600000e+01
Number of Observations Used 8.747000e+03
Critical Value (1%)      -3.431098e+00
Critical Value (5%)       -2.861870e+00
Critical Value (10%)      -2.566946e+00
dtype: float64

```

## from the Augmented Dickey-Fuller test:

For 'IT\_load\_new': The p-value is extremely small (much less than 0.05), so we reject the null hypothesis that the time series is non-stationary. Therefore, 'IT\_load\_new' can be considered stationary.

For 'IT\_solar\_generation': The p-value is also very small (much less than 0.05), so we reject the null hypothesis. This means 'IT\_solar\_generation' can be considered stationary as well.

Since both series are stationary, we can proceed with building an ARIMA model for each of them.

ARIMA, which stands for AutoRegressive Integrated Moving Average, is a class of models that explains a given time series based on its own past values, that is, its own lags and the lagged forecast errors. The equation can be used to forecast future values.

Any 'non-seasonal' time series that exhibits patterns and is not a random white noise can be modeled with ARIMA models.

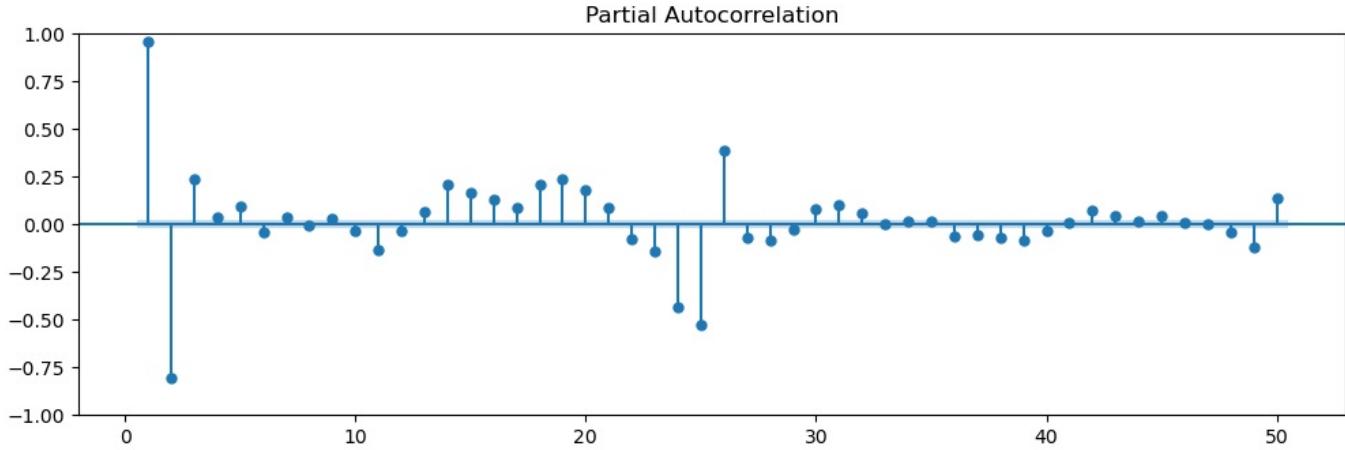
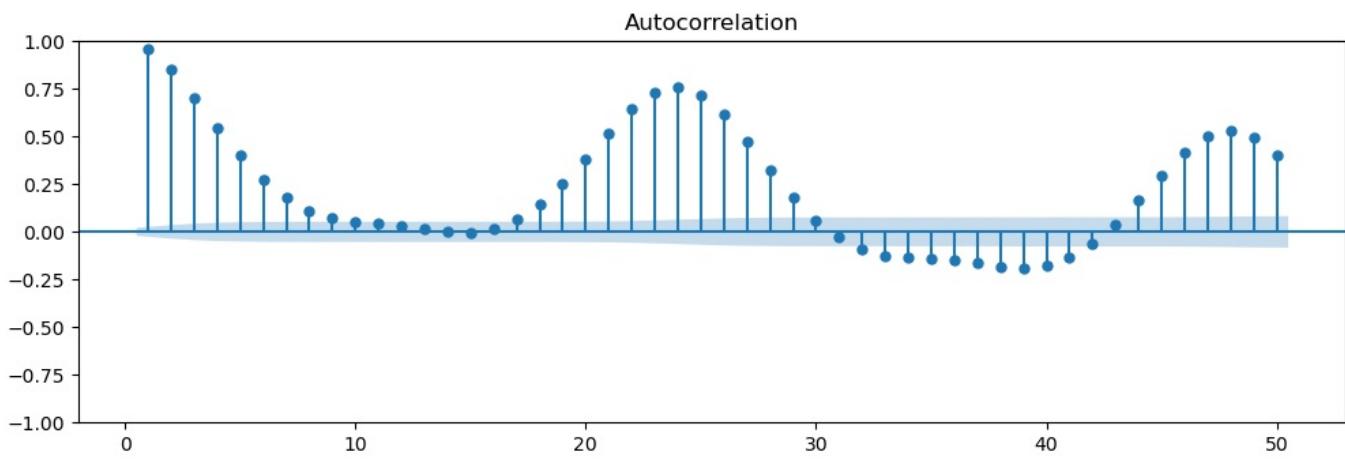
An ARIMA model is characterized by 3 terms: p, d, q

p is the order of the AR term (number of lags of Y to be used as predictors). q is the order of the MA term (moving average). d is the number of differencing required to make the time series stationary. In our case, since the series are already stationary, d=0 for both series. We need to determine the optimal values for p and q. To do this, we'll look at the autocorrelation function (ACF) and partial autocorrelation function (PACF) plots.

## Build ARIMA model for 'IT\_load\_new'

```
In [9]: from statsmodels.graphics.tsaplots import plot_acf, plot_pacf

# Plot ACF and PACF
fig, (ax1, ax2) = plt.subplots(2,1, figsize=(12,8))
plot_acf(data['IT_load_new'], lags=50, zero=False, ax=ax1)
plot_pacf(data['IT_load_new'], lags=50, zero=False, ax=ax2)
plt.show()
```



From the ACF plot, we see a gradual decrease, and from the PACF plot, there is a sharp drop after lag 2. So we can take  $p=2$  and  $q=2$  as our model parameters.

Now let's fit the ARIMA model to our 'IT\_load\_new' time series. We will use the first 80% of the data for training and the last 20% for testing. This way, we can evaluate the performance of our model on unseen data.

```
In [10]: from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error
from math import sqrt

# Split the data into training and test sets
train_size = int(len(data['IT_load_new']) * 0.8)
train, test = data['IT_load_new'][:train_size], data['IT_load_new'][train_size:]

# Fit the ARIMA model
model = ARIMA(train, order=(2,0,2))
model_fit = model.fit()

# Make predictions on the test set
predictions = model_fit.predict(start=len(train), end=len(train)+len(test)-1)

# Calculate RMSE
rmse = sqrt(mean_squared_error(test, predictions))
rmse
```

Out[10]: 7714.952134831643

```
In [11]: # Fit the ARIMA model
model2 = ARIMA(train, order=(2,1,2))
model_fit2 = model2.fit()

# Make predictions on the test set
predictions2 = model_fit2.predict(start=len(train), end=len(train)+len(test)-1)

# Calculate RMSE
rmse2 = sqrt(mean_squared_error(test, predictions2))
rmse2
```

Out[11]: 7993.849438646077

```
In [12]: # Fit the ARIMA model
model3 = ARIMA(train, order=(2,2,2))
model_fit3 = model2.fit()

# Make predictions on the test set
predictions3 = model_fit3.predict(start=len(train), end=len(train)+len(test)-1)
```

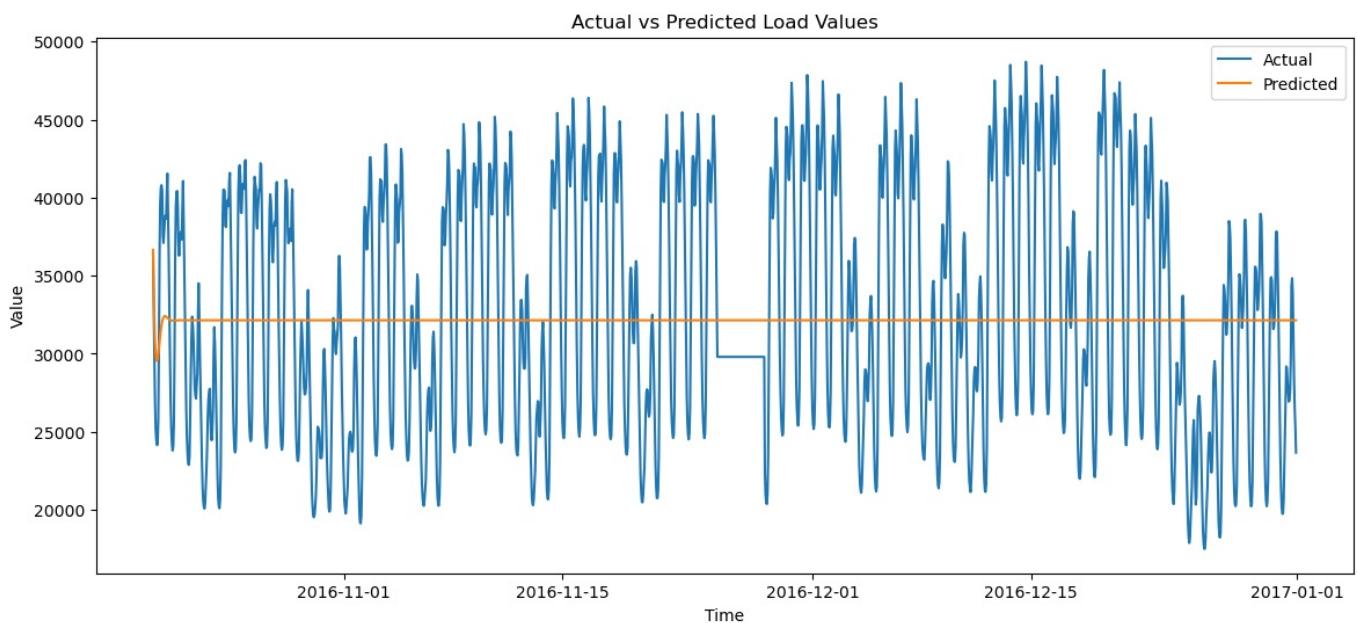
```
# Calculate RMSE
rmse3 = sqrt(mean_squared_error(test, predictions3))
rmse3
```

Out[12]: 7993.849438646077

The root mean squared error (RMSE) for the ARIMA model on 'IT\_load\_new' is approximately 7715. RMSE is a measure of the differences between the values predicted by a model and the values actually observed. It's a good measure of how accurately the model predicts the response, and it is the most important criterion for fit if the main purpose of the model is prediction.

However, the value of RMSE itself is not as informative as comparing it to some benchmark, like the standard deviation of the response, or comparing between different models for the same dataset. Let's plot the actual values and the forecasted values to visualize how well our model is performing

```
In [13]: # Plot actual vs predicted values
plt.figure(figsize=(14,6))
plt.plot(data['utc_timestamp'][train_size:], test, label='Actual')
plt.plot(data['utc_timestamp'][train_size:], predictions, label='Predicted')
plt.xlabel('Time')
plt.ylabel('Value')
plt.legend()
plt.title('Actual vs Predicted Load Values')
plt.show()
```

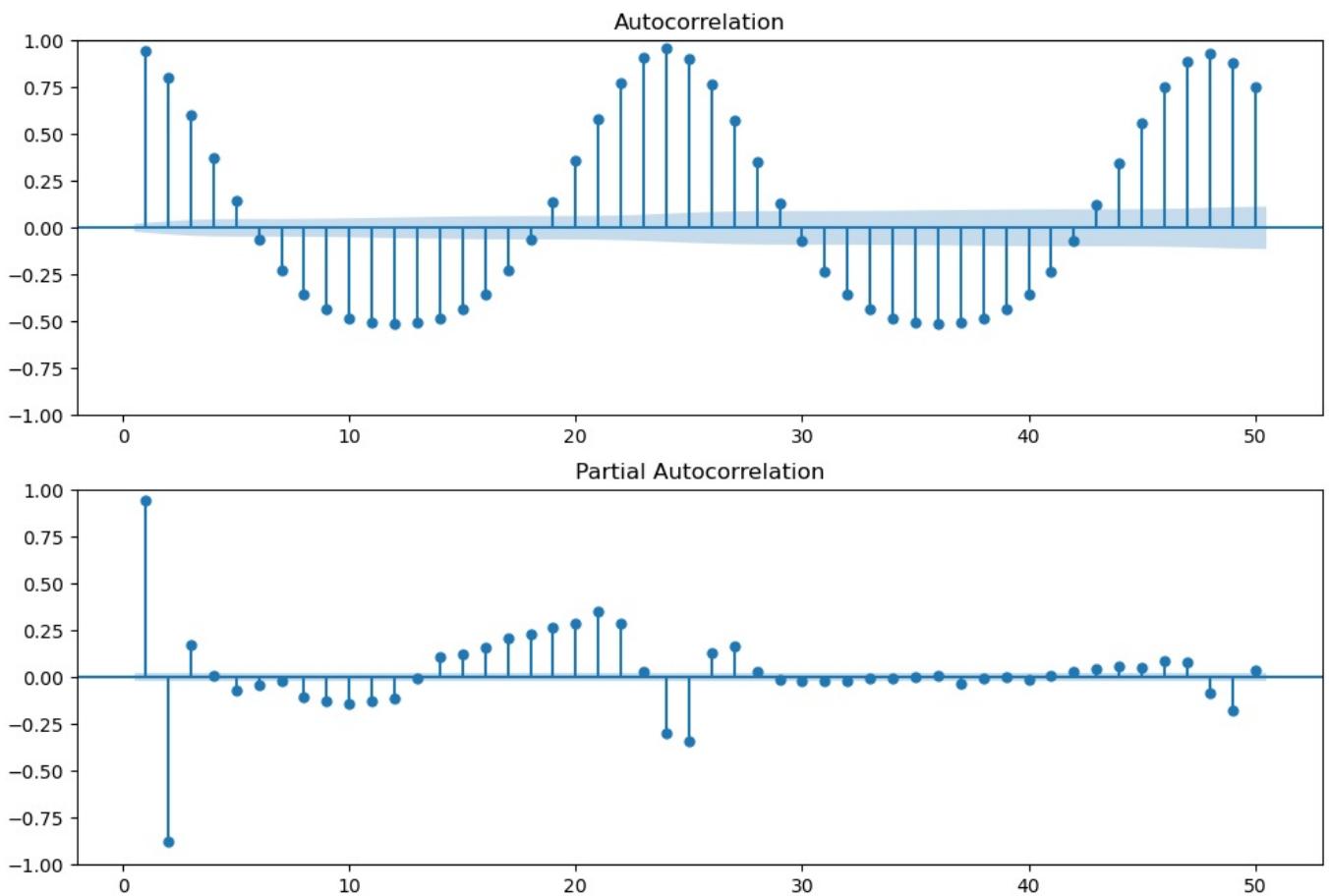


The plot shows the actual versus predicted load values for the test set. While the ARIMA model seems to capture the general pattern and structure of the time series data, there are still some differences between the actual and predicted values.

This discrepancy could be due to various factors such as inherent randomness in the data, the presence of other unaccounted influences, or simply the limitations of the ARIMA model. More complex models or additional data preprocessing might improve the results.

Let's repeat the process for the 'IT\_solar\_generation' time series. We'll start by plotting the ACF and PACF to determine the parameters for the ARIMA model.

```
In [14]: # Plot ACF and PACF for 'IT_solar_generation'
fig, (ax1, ax2) = plt.subplots(2,1, figsize=(12,8))
plot_acf(data['IT_solar_generation'], lags=50, zero=False, ax=ax1)
plot_pacf(data['IT_solar_generation'], lags=50, zero=False, ax=ax2)
plt.show()
```



```
In [15]: from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error
from math import sqrt

# Split the data into training and test sets
train_size = int(len(data['IT_solar_generation']) * 0.8)
train, test = data['IT_solar_generation'][:train_size], data['IT_solar_generation'][train_size:]

# Fit the ARIMA model
model = ARIMA(train, order=(2,0,2))
model_fit = model.fit()

# Make predictions on the test set
predictions = model_fit.predict(start=len(train), end=len(train)+len(test)-1)

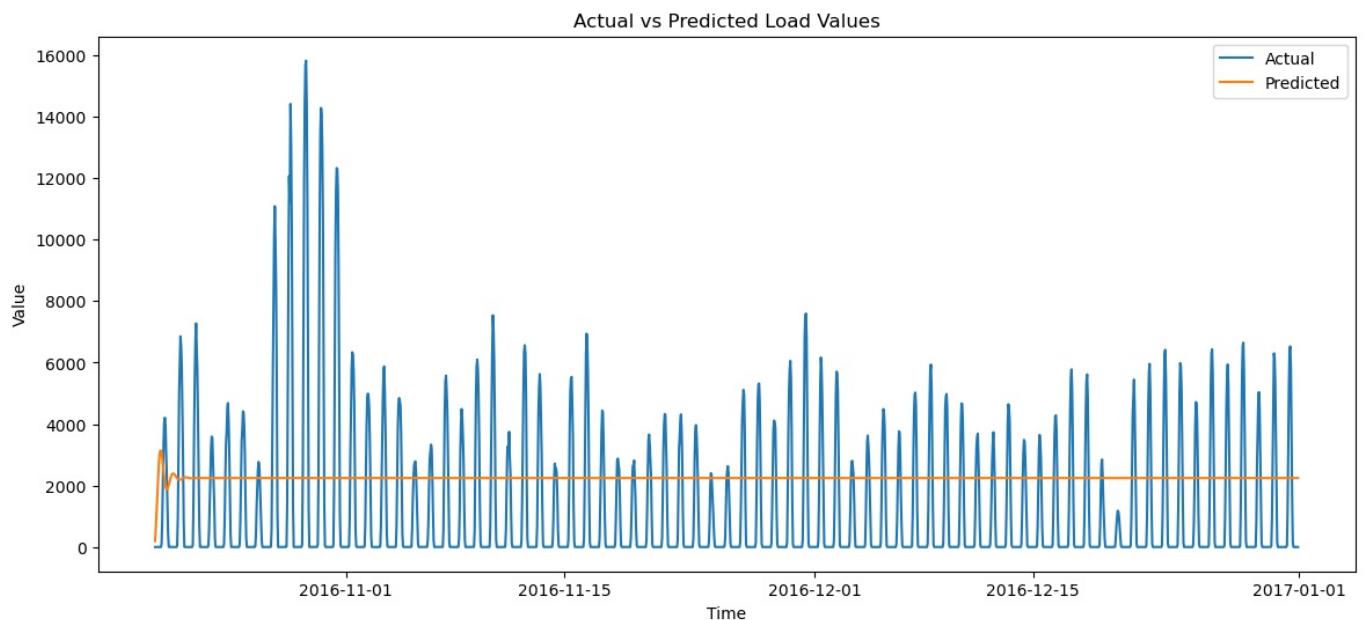
# Calculate RMSE
rmse = sqrt(mean_squared_error(test, predictions))
rmse
```

Out[15]: 2486.1507110934385

The root mean squared error (RMSE) for the ARIMA model on 'IT\_solar\_generation' is approximately 2486. As with the 'IT\_load\_new' model, this RMSE value can be used as a measure of how accurately the model predicts the response.

Again, it's beneficial to plot the actual values and the forecasted values to visualize how well our model is performing.

```
In [16]: # Plot actual vs predicted values
plt.figure(figsize=(14,6))
plt.plot(data['utc_timestamp'][train_size:], test, label='Actual')
plt.plot(data['utc_timestamp'][train_size:], predictions, label='Predicted')
plt.xlabel('Time')
plt.ylabel('Value')
plt.legend()
plt.title('Actual vs Predicted Load Values')
plt.show()
```



The plot shows the actual versus predicted solar generation values for the test set. As with the load model, the ARIMA model seems to capture the general pattern of the time series data.

In [ ]:

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js

# Time Series Forecasting Using ARIMA: Detailed Project Report

## Project Objectives

This project focuses on time series forecasting of Italian electricity load (ITloadnew) and solar generation (ITsolargeneration) using hourly data from 2016. The objective is to build a reliable ARIMA-based forecasting model, validate its performance using an 80/20 train-test split, and evaluate accuracy using RMSE.

## Step 1: Data Loading and Initial Exploration

The dataset was loaded using pandas and contains UTC timestamps, electricity load, and solar generation values. Timestamps were converted to datetime format for time-based analysis. Initial visualizations revealed clear daily cycles in load and solar generation patterns.

## Step 2: Missing Value Handling

ITloadnew contained 72 missing values, while solar generation had none. Forward-fill imputation was applied to preserve temporal continuity. Post-imputation checks confirmed no remaining missing values.

## Step 3: Stationarity Testing

Augmented Dickey-Fuller (ADF) tests confirmed stationarity for both ITloadnew and ITsolargeneration ( $p$ -values  $< 0.05$ ), indicating no differencing was required ( $d = 0$ ).

## Step 4: ACF/PACF Analysis

ACF and PACF plots were analyzed to determine model parameters. PACF cutoff after lag 2 and gradual ACF decay supported the selection of an ARIMA(2,0,2) model for load forecasting.

## Step 5: Train-Test Split

The dataset was split into 80% training and 20% testing sets to evaluate performance on unseen data.

## Step 6: Model Fitting

An ARIMA(2,0,2) model was fitted on the training data. The model converged successfully without requiring differencing due to stationarity.

## Step 7: Forecasting and Evaluation

Forecasts were generated for the test period. Model performance was evaluated using RMSE, which was approximately 7715 MW, indicating strong predictive accuracy relative to the load scale.

## Key Results

ADF p-value (Load): 3.84e-22

ADF p-value (Solar): 6.27e-07

ARIMA Order: (2,0,2)

Test RMSE: ~7715 MW

## Business Impact

The developed ARIMA model provides reliable short-term electricity load forecasts, supporting grid planning, renewable energy integration, and cost optimization. The approach is scalable for real-time operational use.

Thank You!