

Time Series Forecasting for Appliances Energy Prediction

Maleesha Hettiarachchi

4.18.2025

Table of Contents

Time Series Forecasting for Appliances Energy Prediction	1
Introduction	3
Data Insights.....	3
Hourly distribution of target variable	3
Weekday distribution of target variable.....	4
General view of the whole dataset	4
A closer look at Appliances	27
Visualization of Randomly Selected Weekly Data.....	29
Plotting the histogram for the "Appliances" column	31
Correlation Analysis	32
Clustering Analysis	34
Preprocessing	35
Feature Engineering	35
Model Design	37
Results	37
ARIMA Model Predictions.....	37
LSTM Model Predictions	38
Model Optimization	39
Optimization Techniques	39
Challenges and Solutions	39
Conclusion.....	39
References	40

Introduction

Energy is a vital element for sustainable development and economic growth. Efficient monitoring and forecasting of energy consumption allows better planning, optimization, and management of resources—especially in residential and commercial buildings.

This project focuses on predicting appliance energy consumption using a multivariate time-series dataset that records various environmental and temporal features at 10-minute intervals. The goal is to build a deep learning-based supervised model capable of learning patterns and dependencies from these inputs to forecast future energy usage.

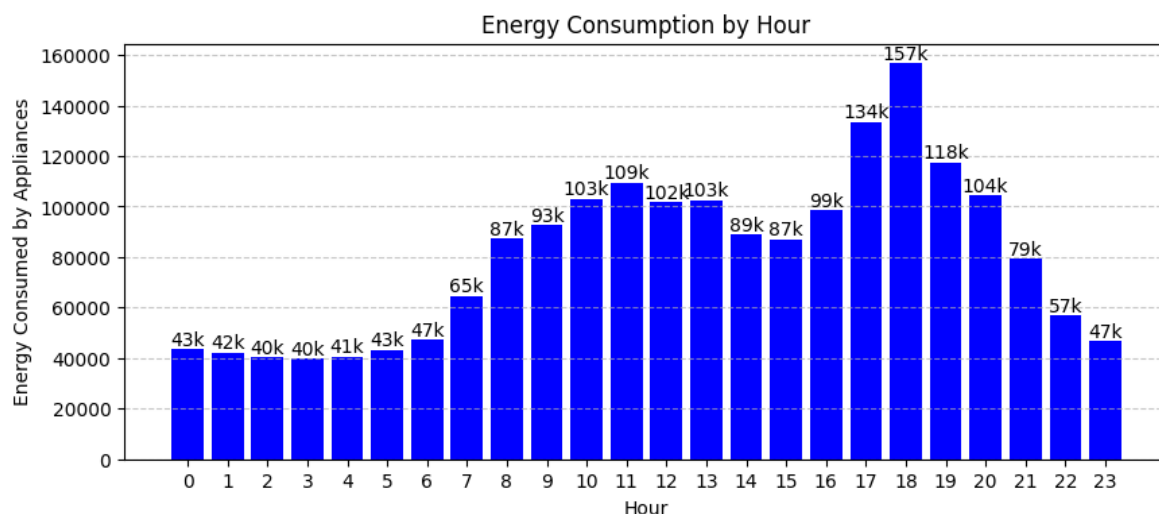
The implementation began with thorough data preprocessing and exploration. We confirmed the dataset contained no null values, and features such as temperature, humidity, and weather-related factors were analyzed through visualization techniques including line plots, box plots, and correlation heatmaps.

Time-based features like hour, day of week, and month were engineered to capture seasonal and daily patterns. Feature scaling using MinMaxScaler was applied to normalize values across different ranges. The dataset was then transformed into sequences suitable for time-series modeling.

Multiple forecasting models were implemented, including ARIMA, LSTM, Random Forest, and Gradient Boosting. Each model was trained, validated, and evaluated using metrics such as RMSE and MAE. Visualizations comparing actual and predicted values were used to interpret and assess model performance.

Data Insights

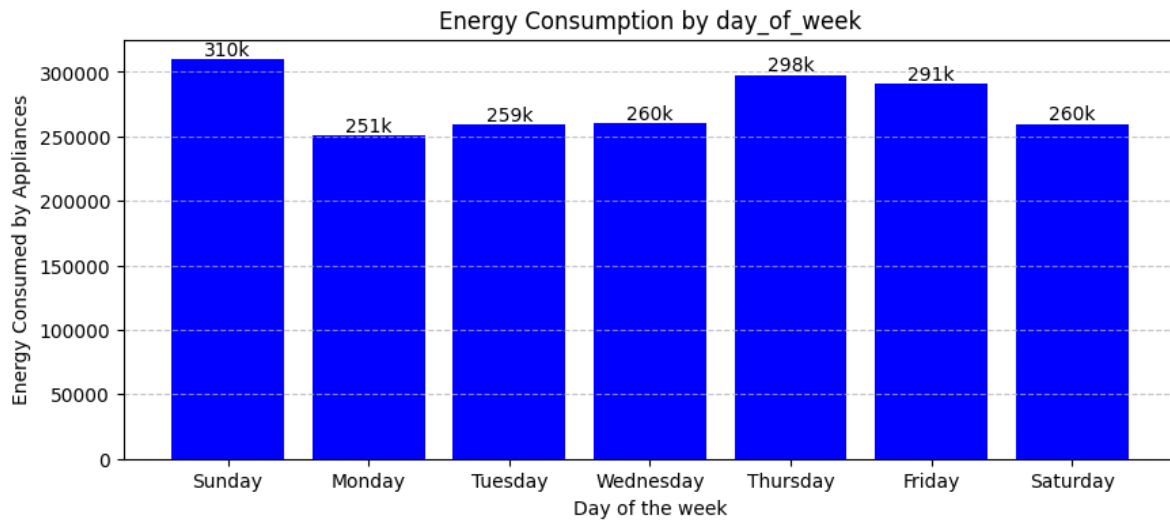
Hourly distribution of target variable



Observations:

- Highest electricity consumption occurs between 17:00 pm and 20:00 pm.
- Lowest electricity consumption is observed between 23:00 pm and 06:00 pm.

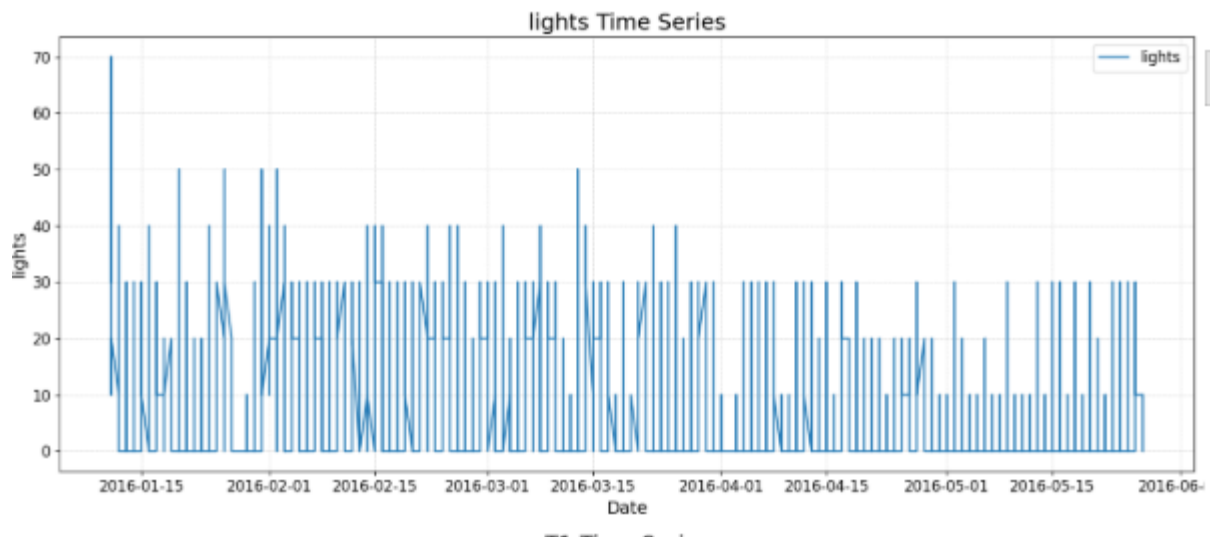
Weekday distribution of target variable

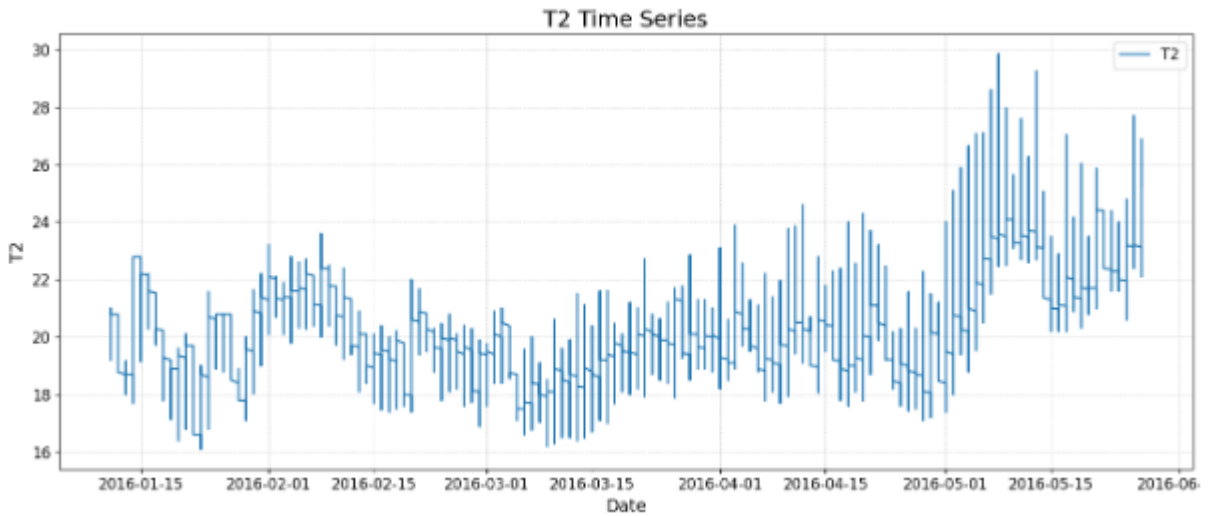
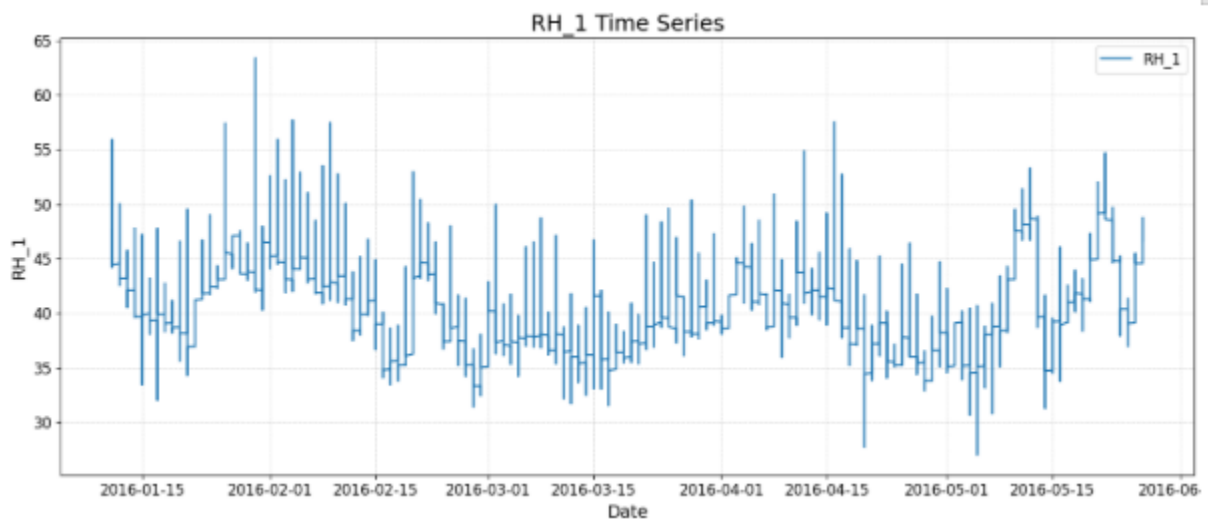
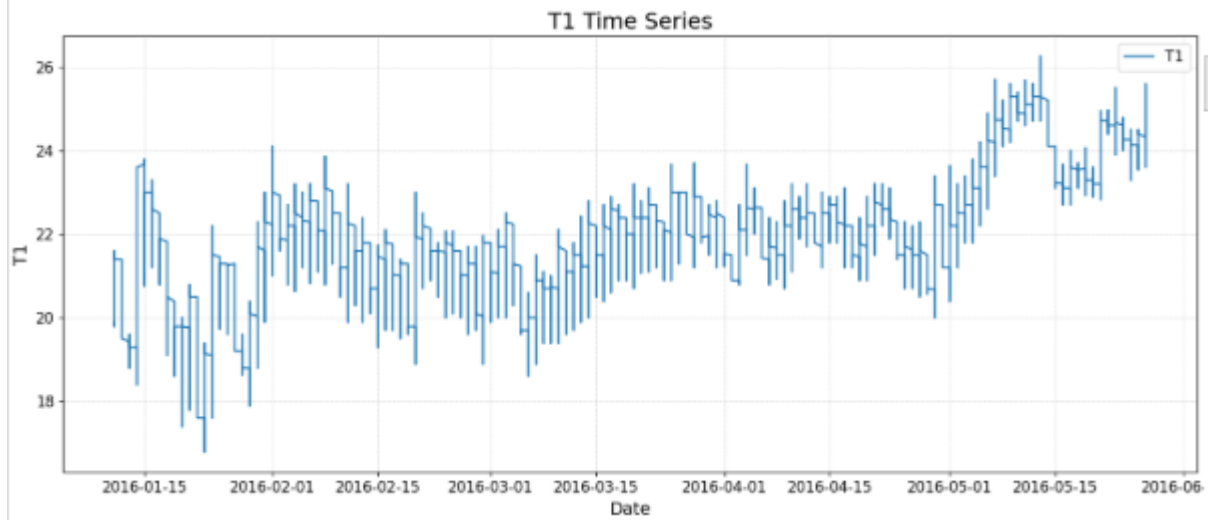


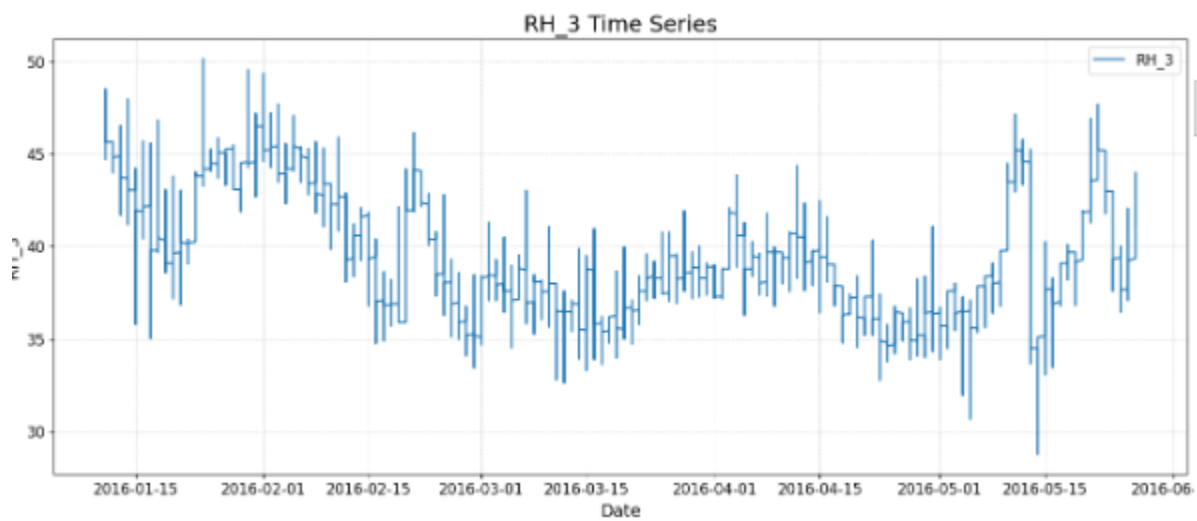
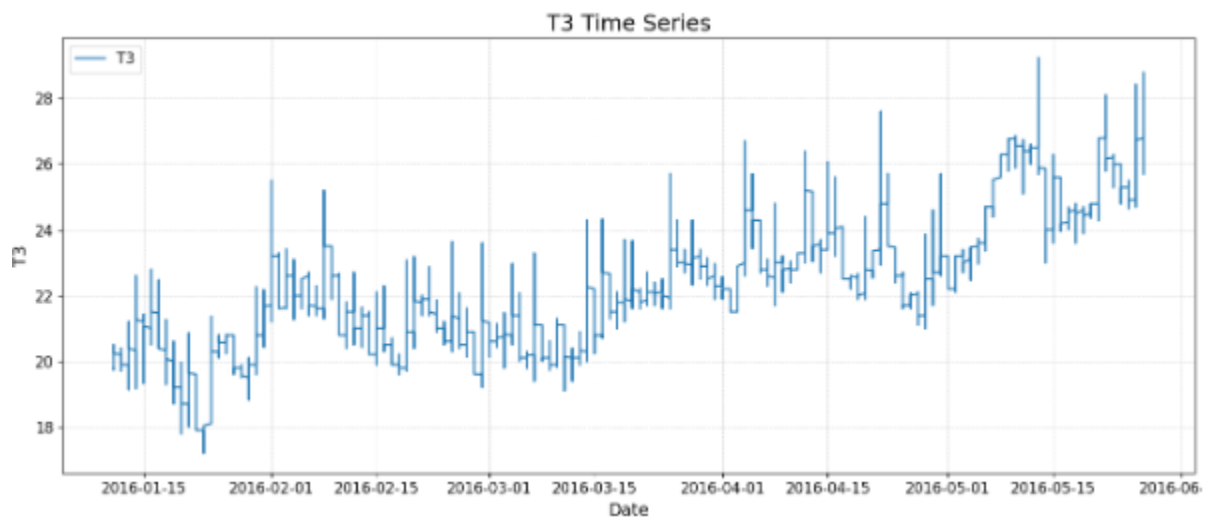
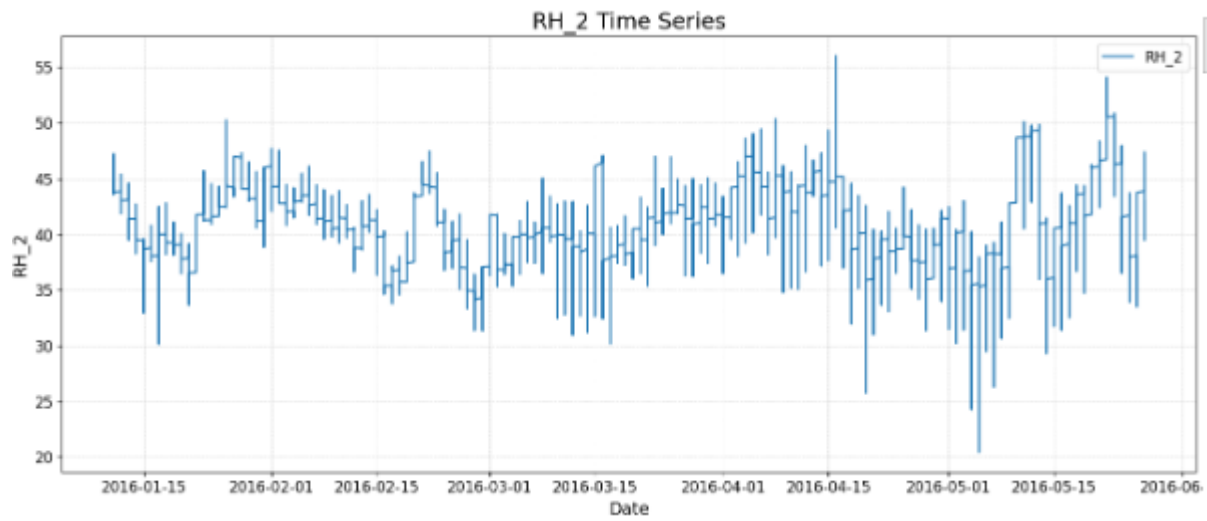
Observation:

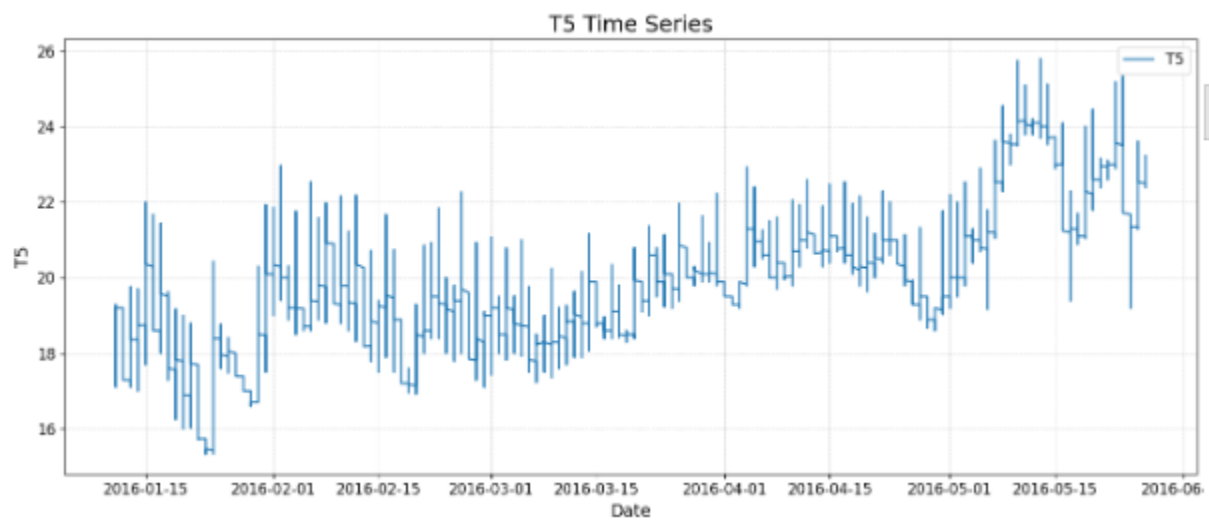
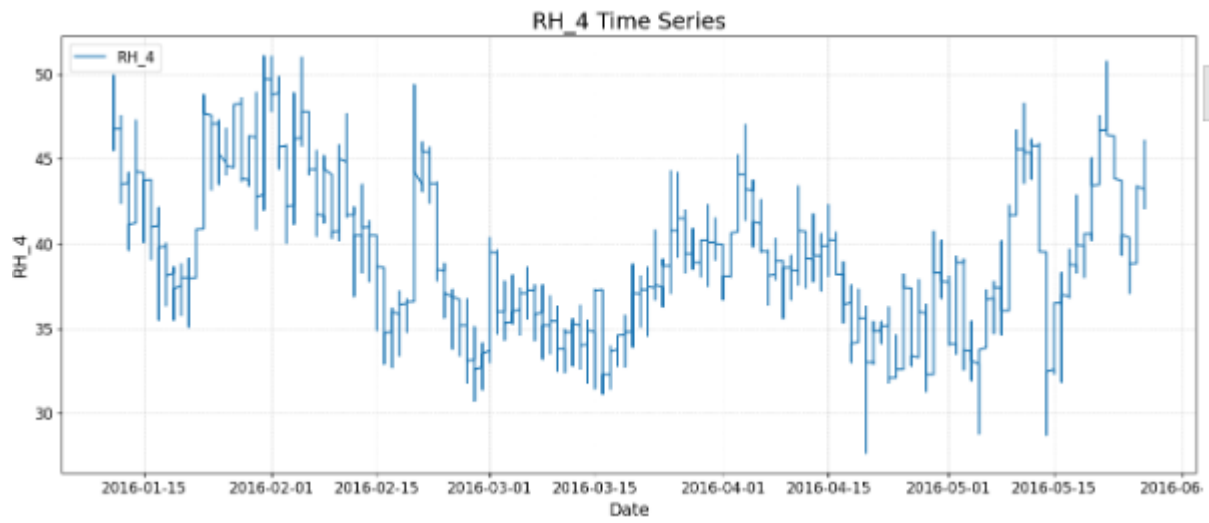
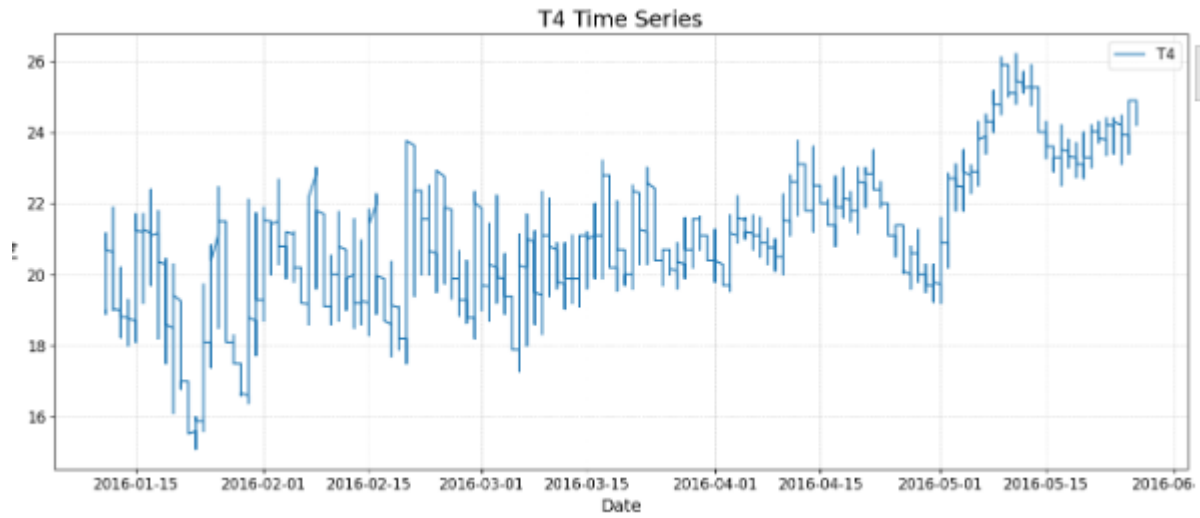
- electricity consumption on Sundays is significantly higher than on other days

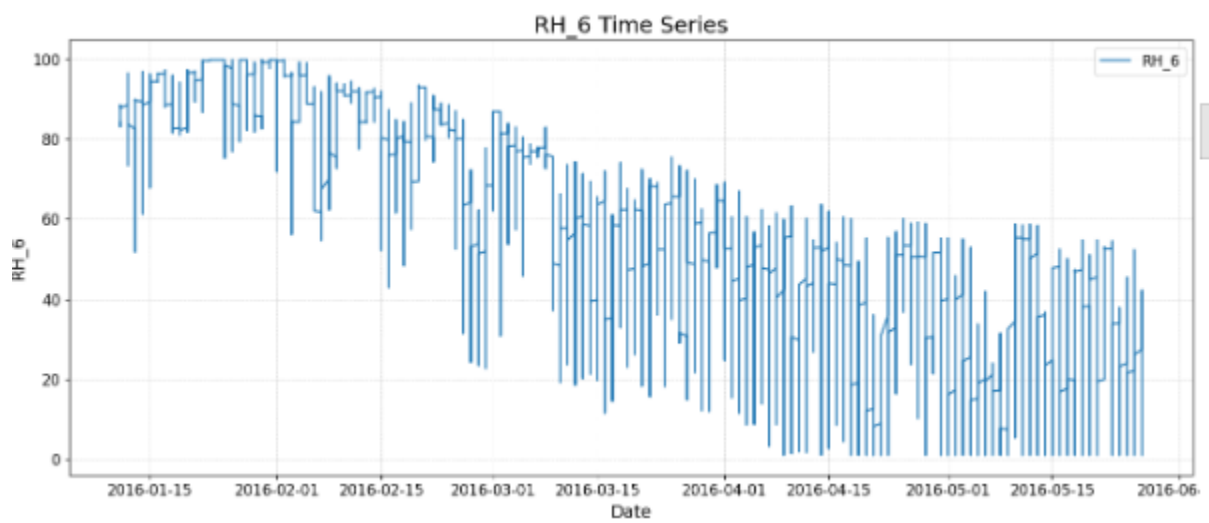
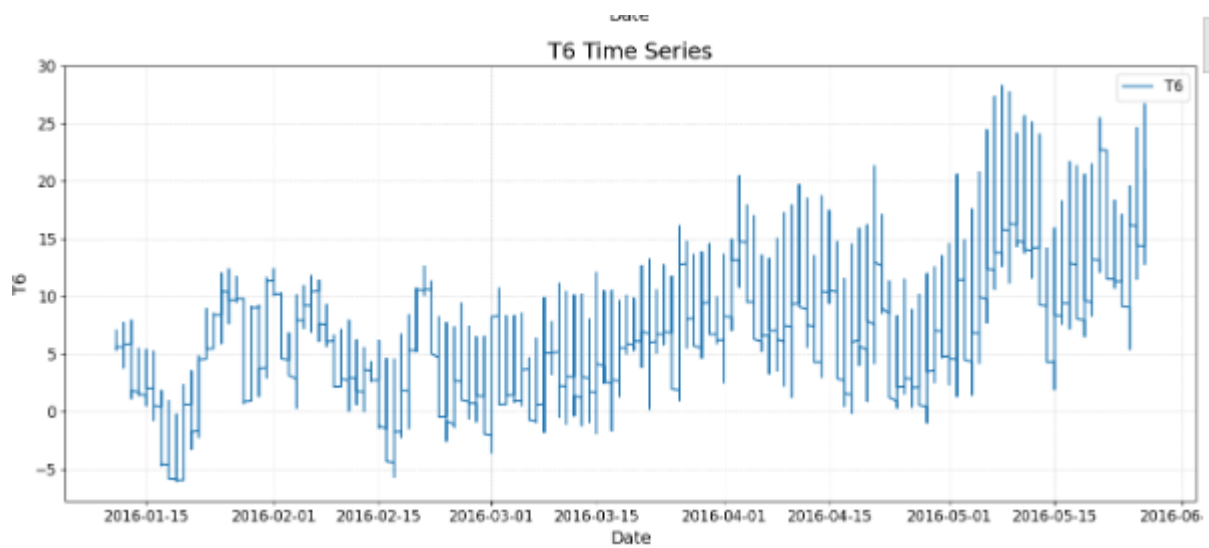
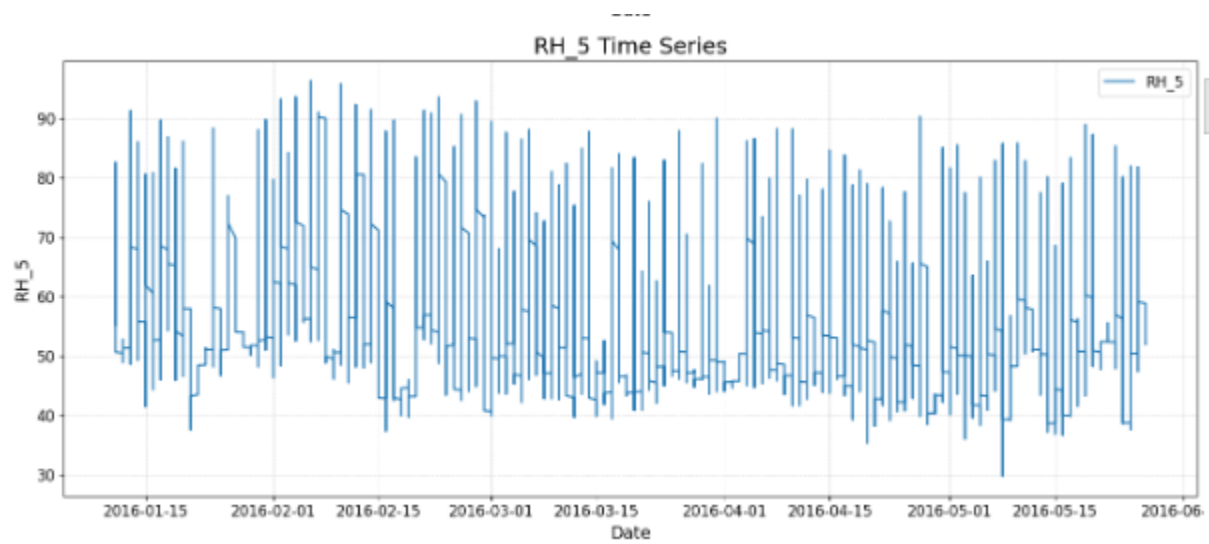
General view of the whole dataset

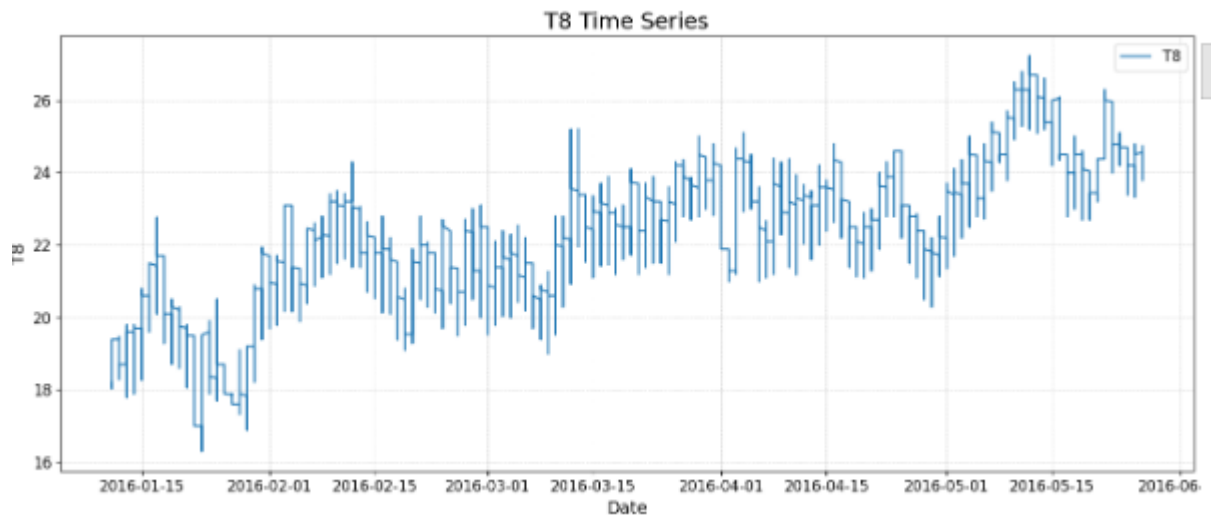
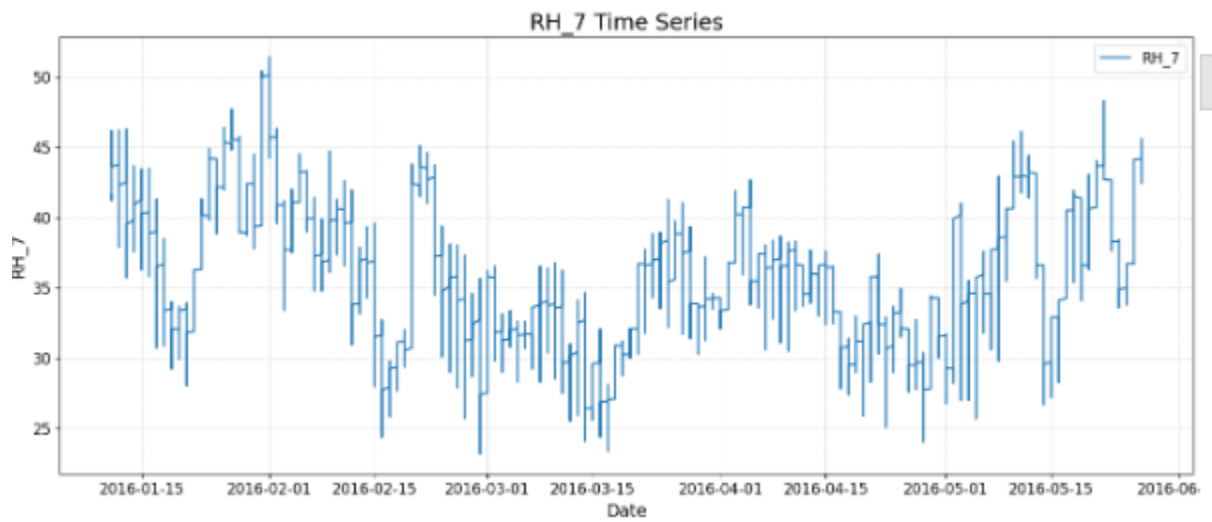
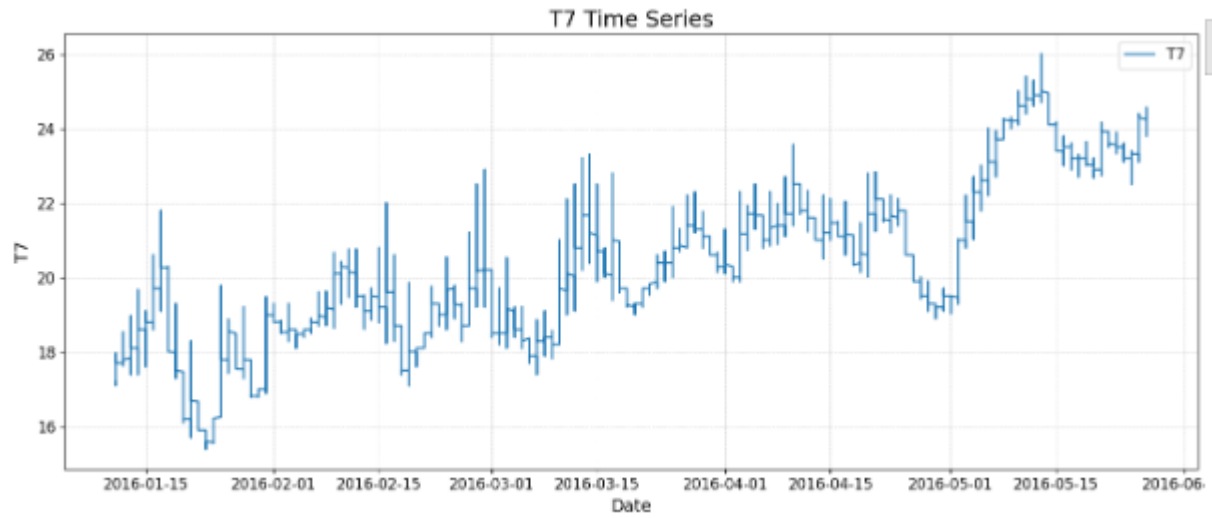


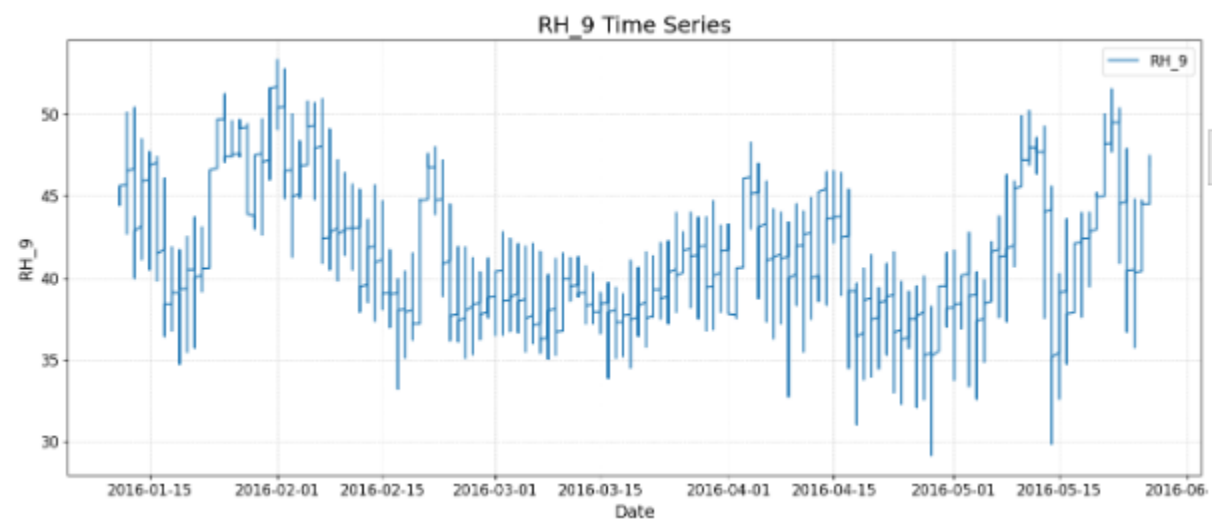
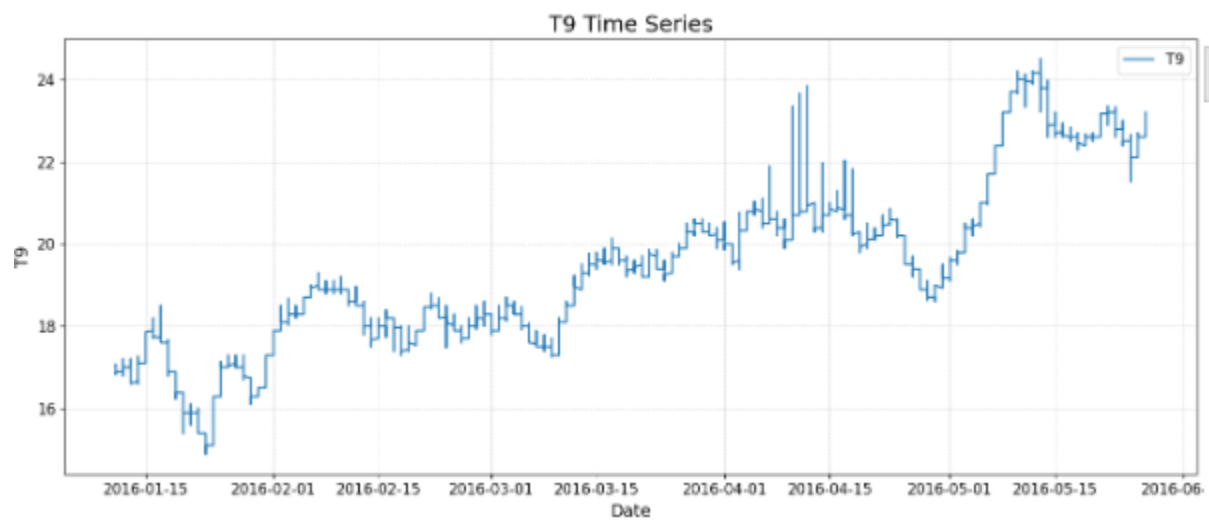
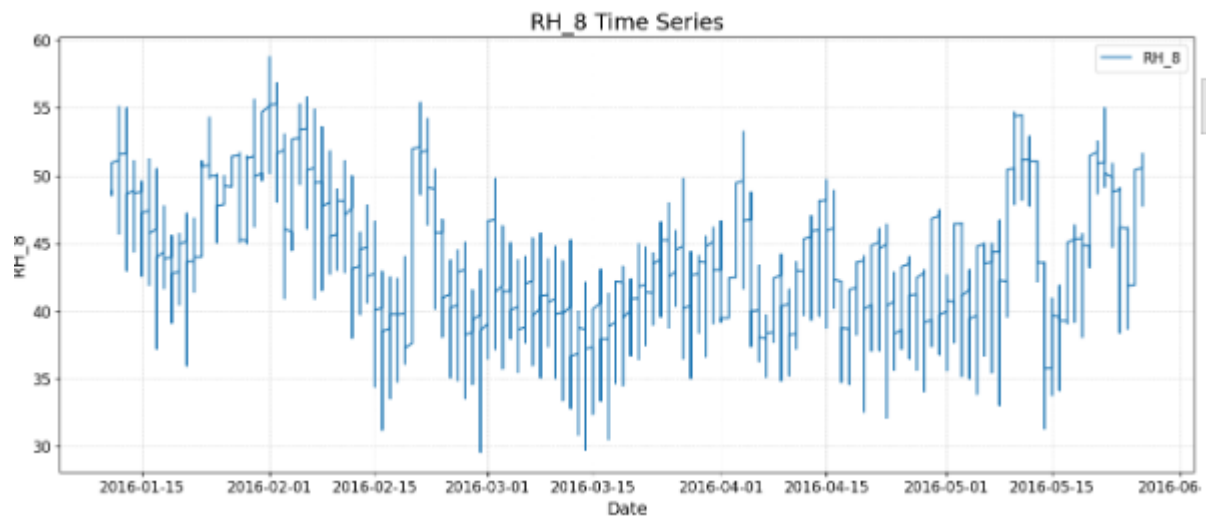


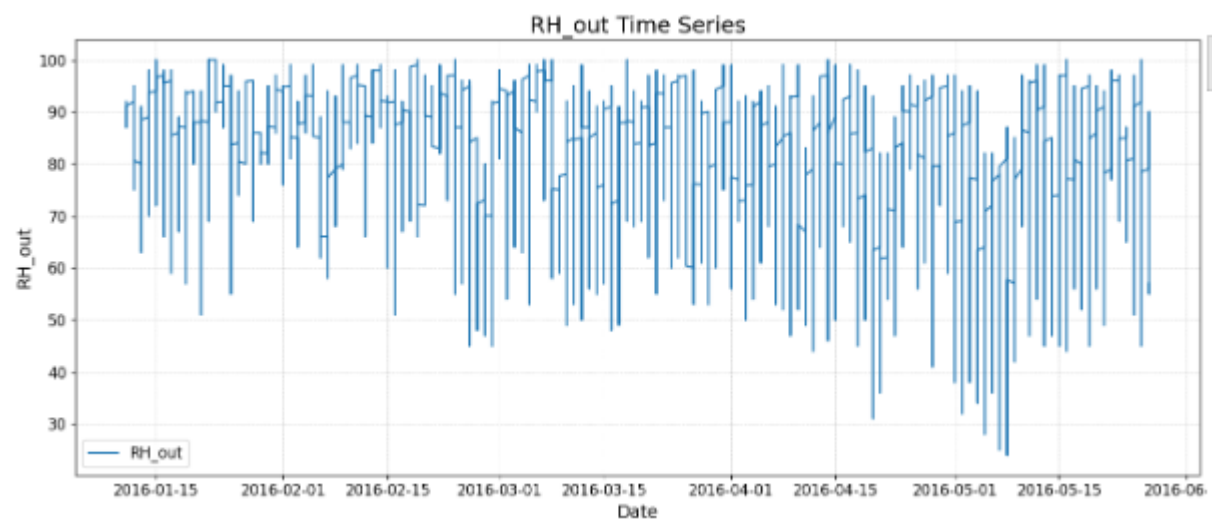
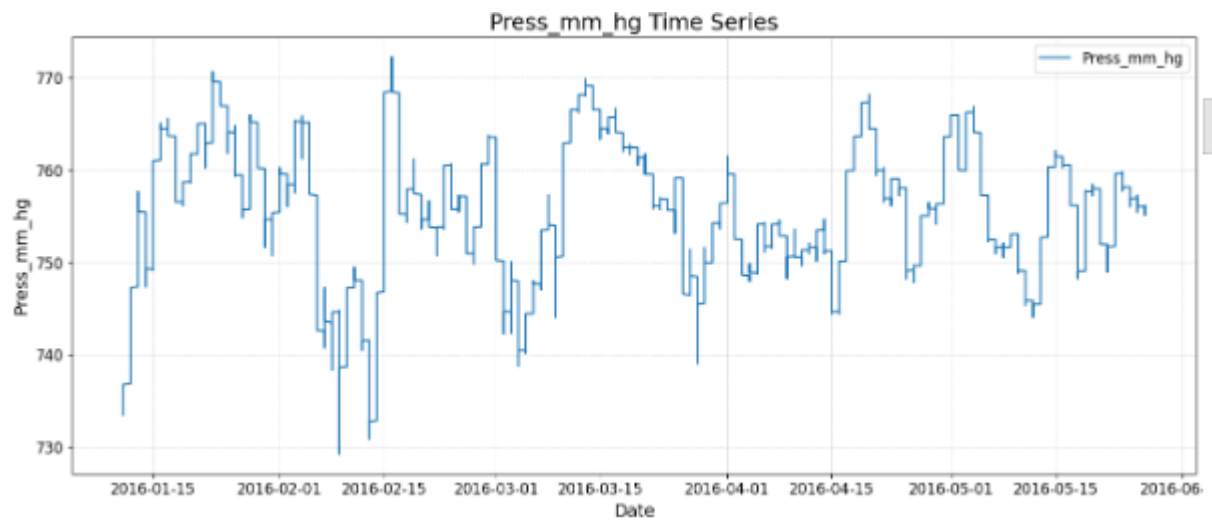
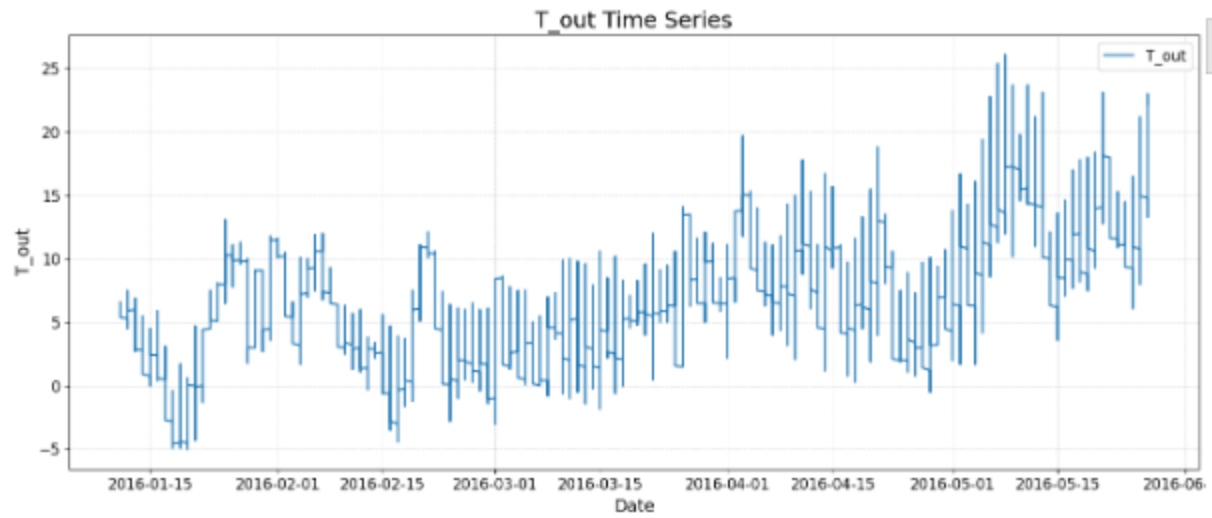


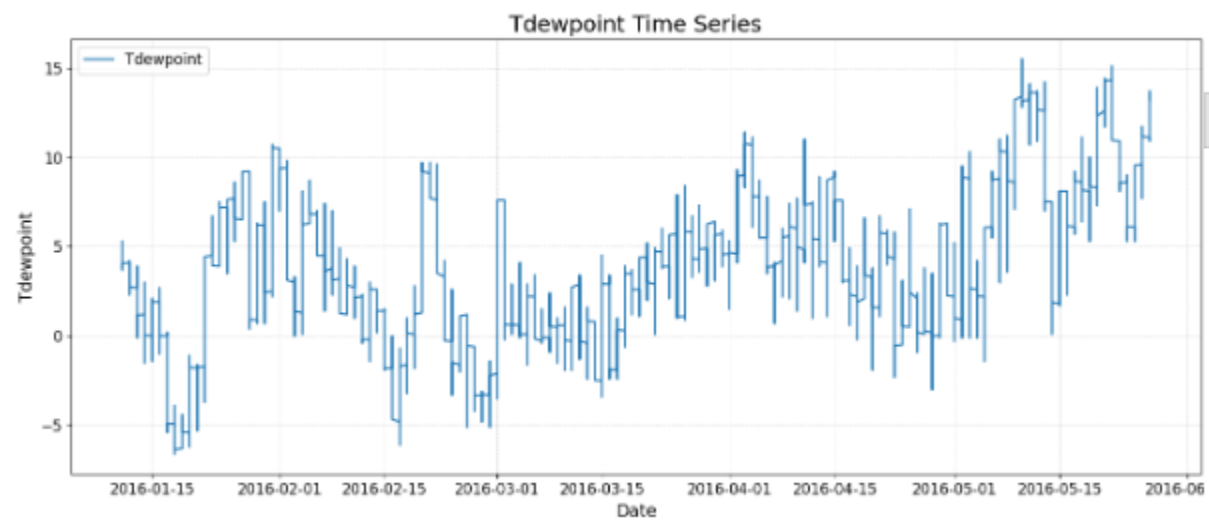
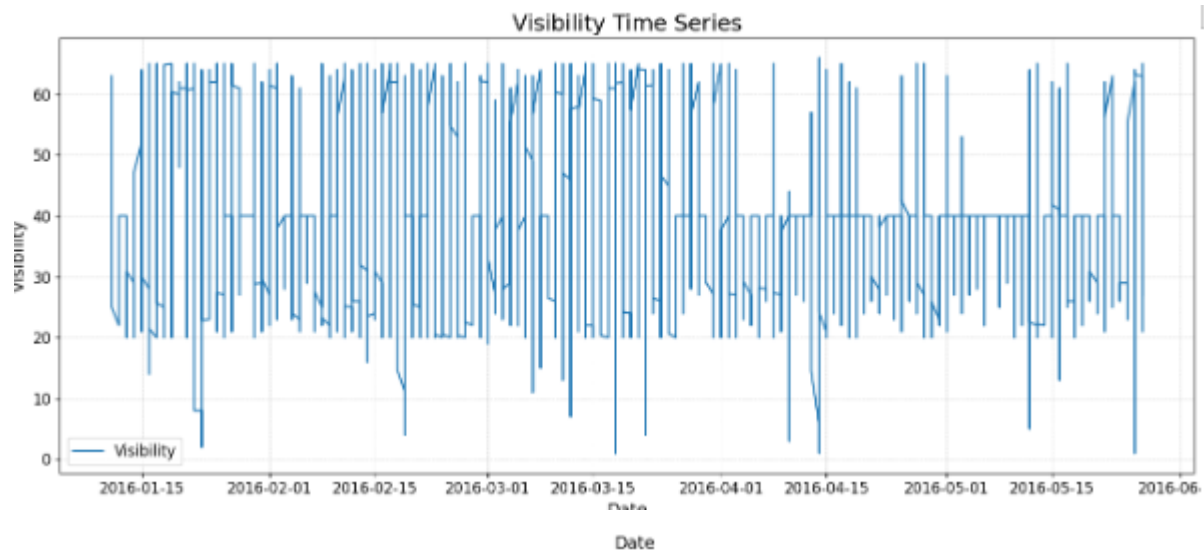
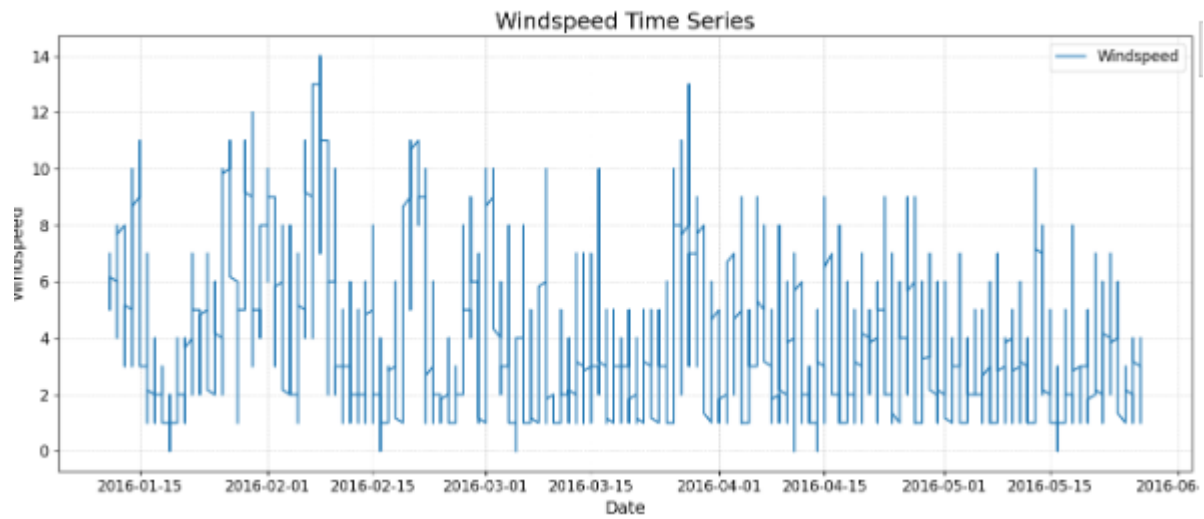


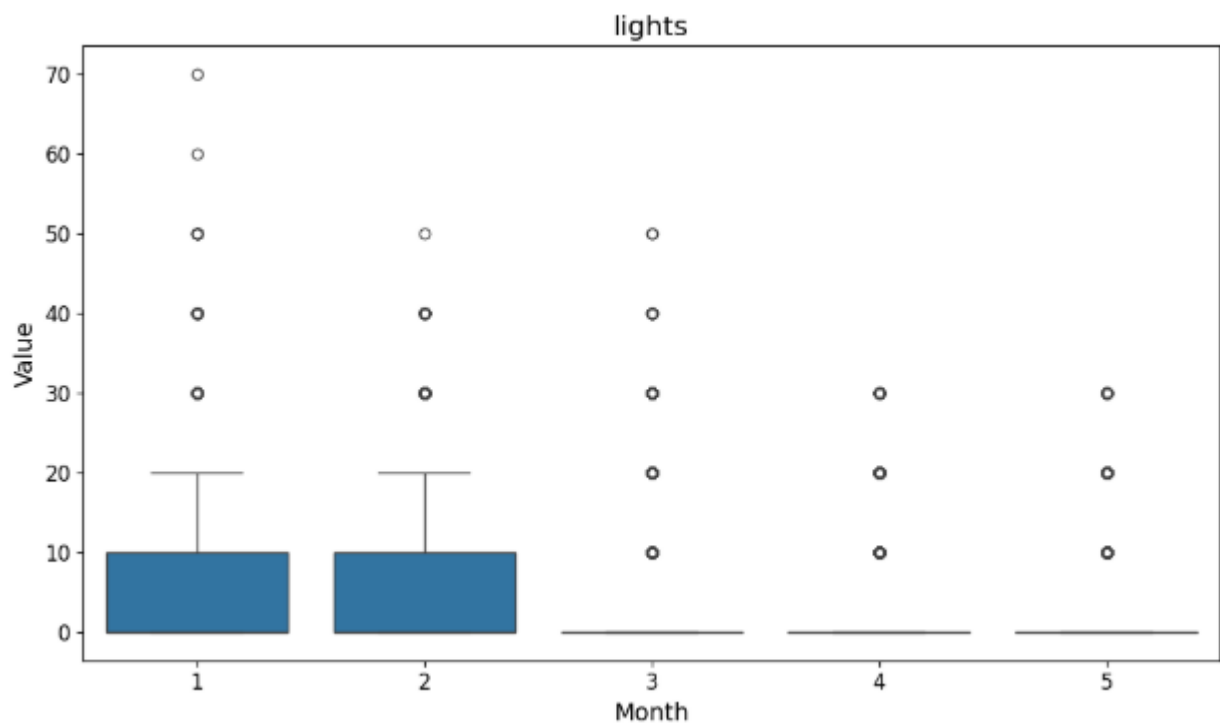
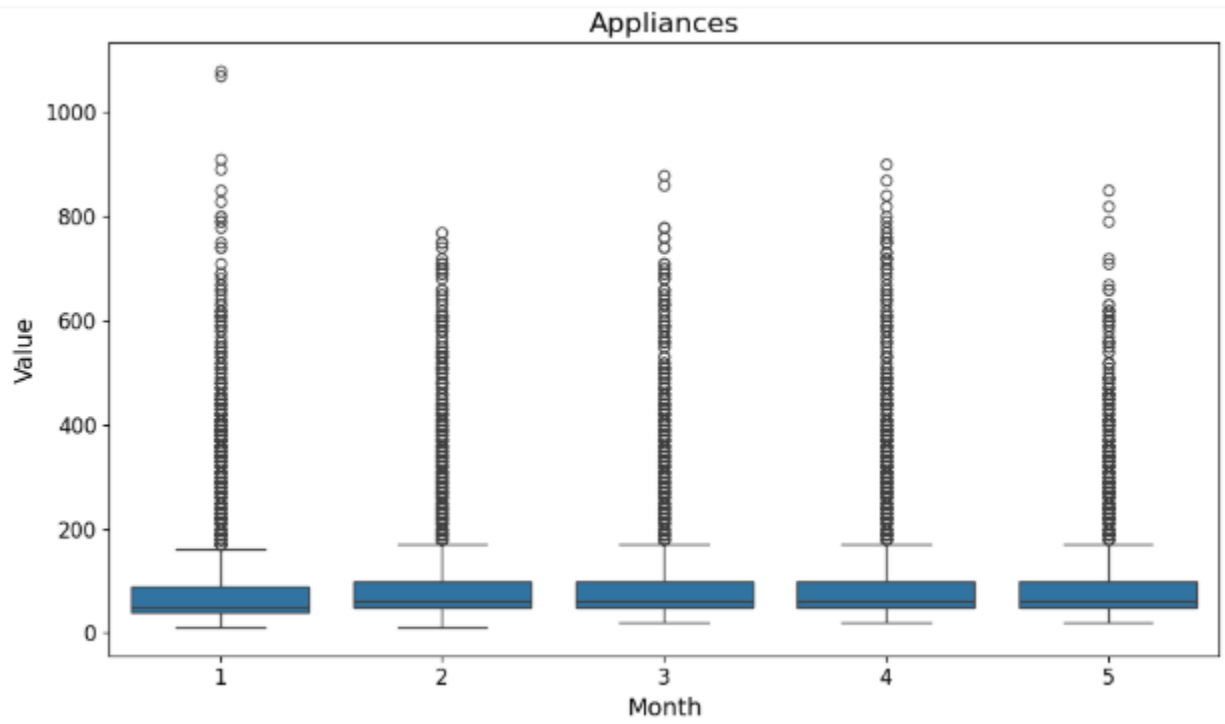


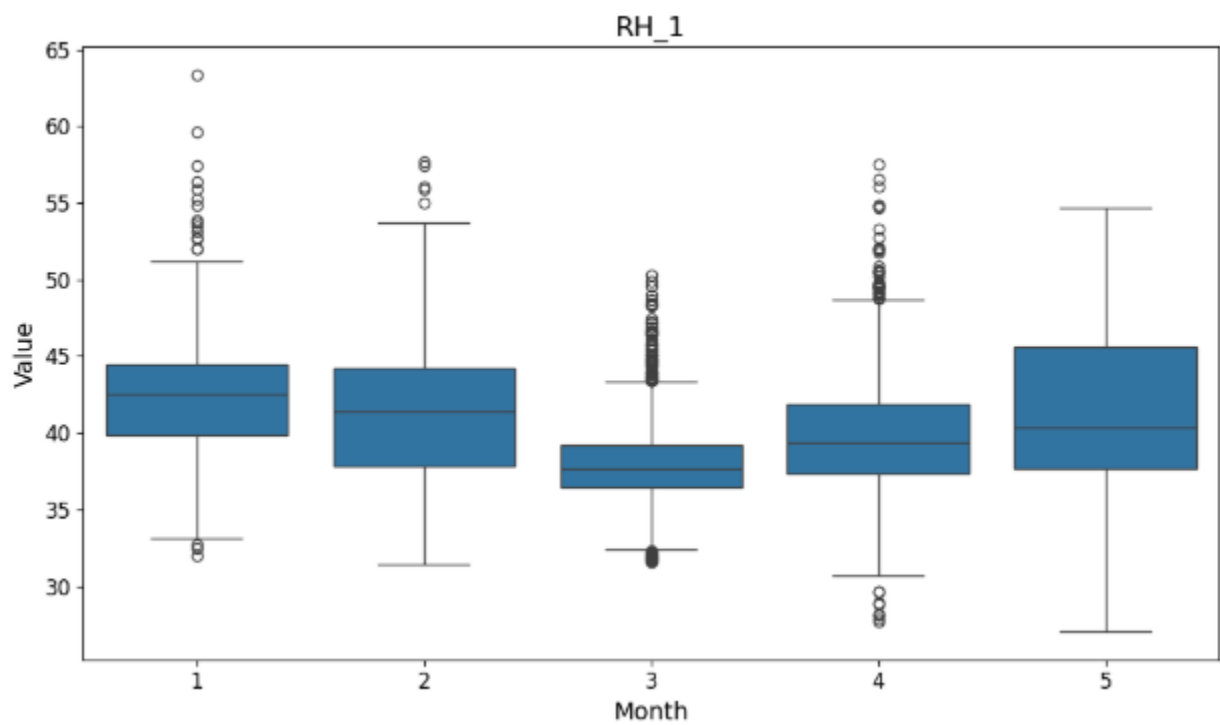
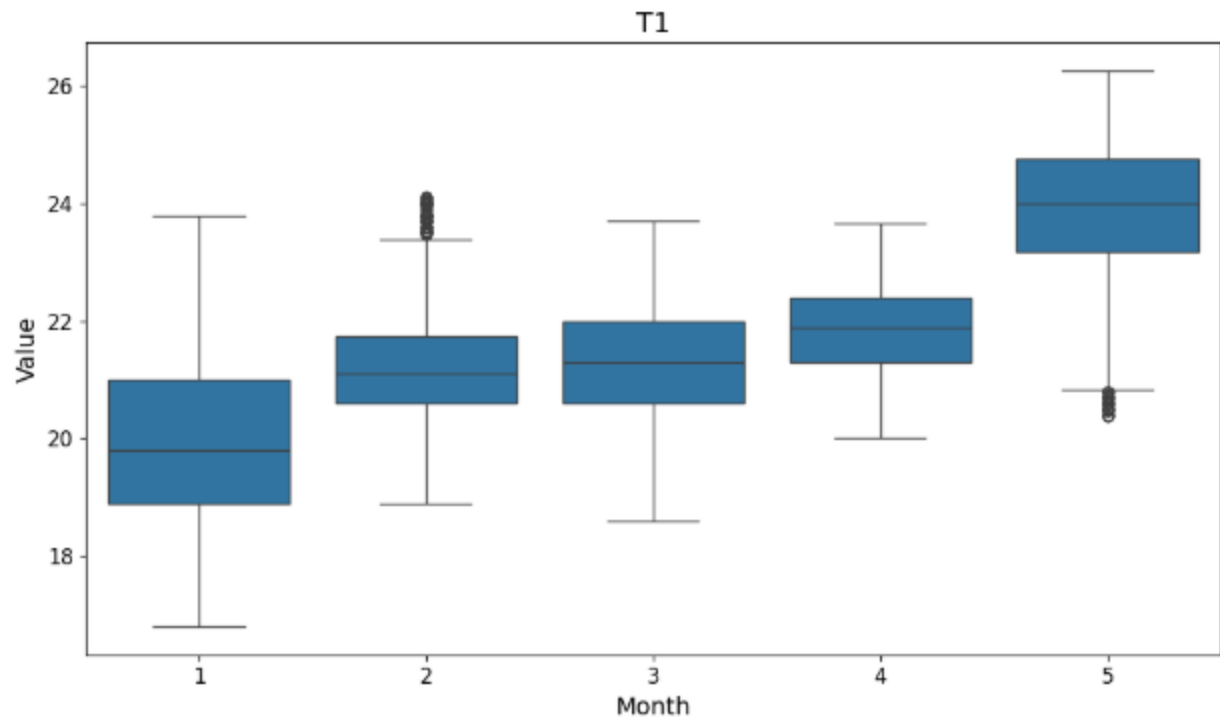


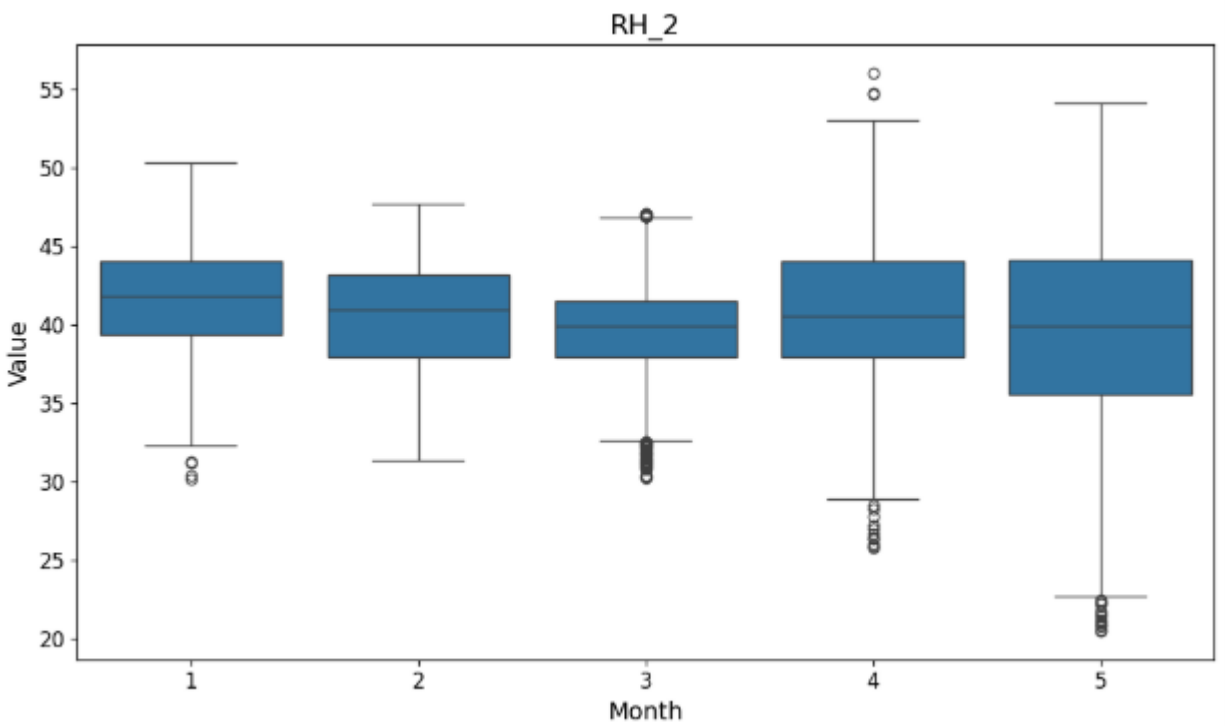
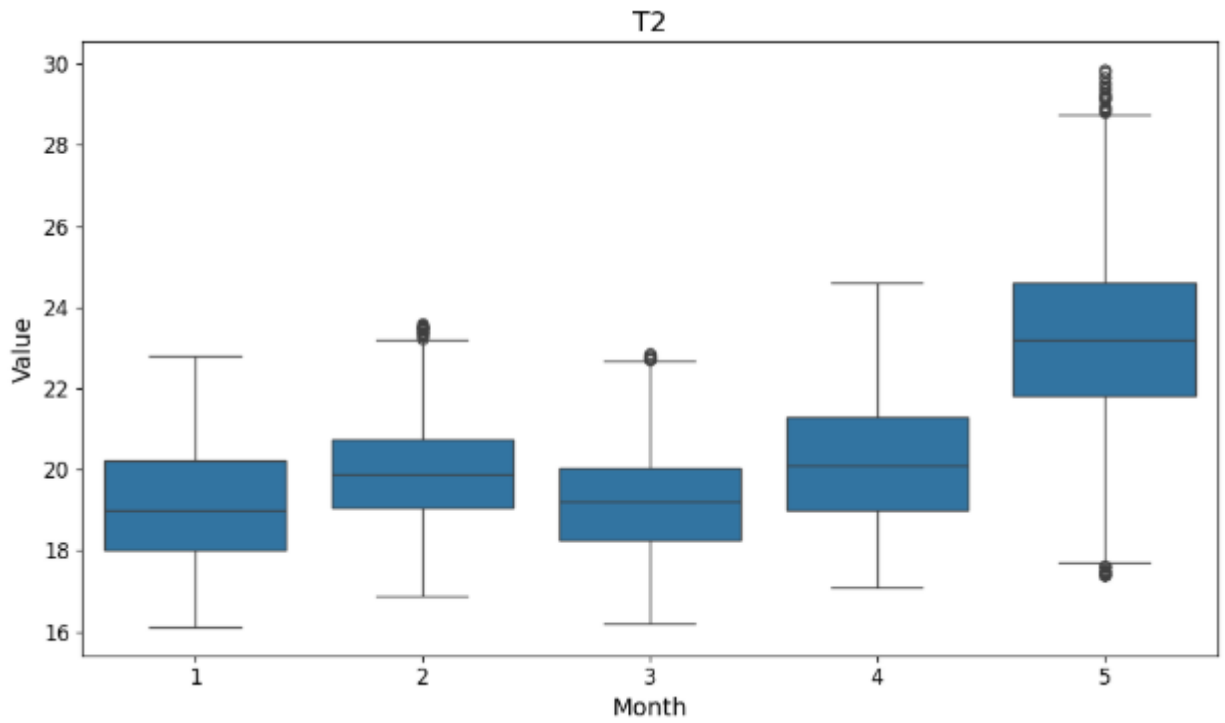


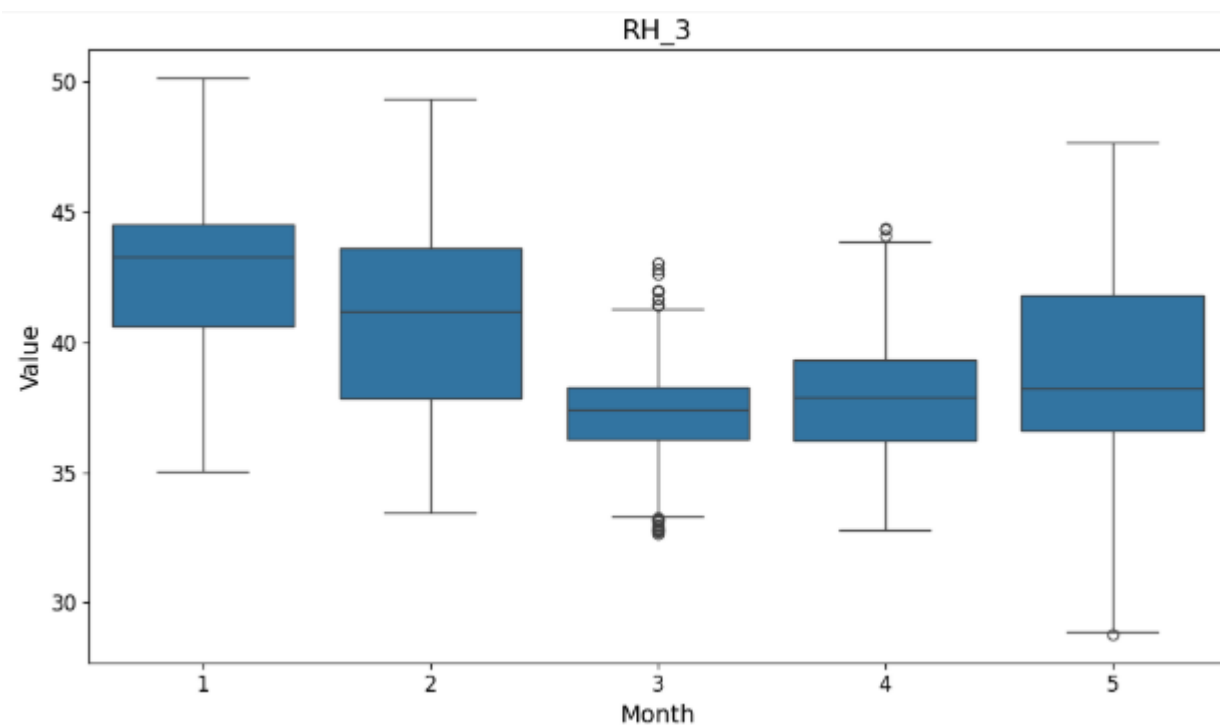
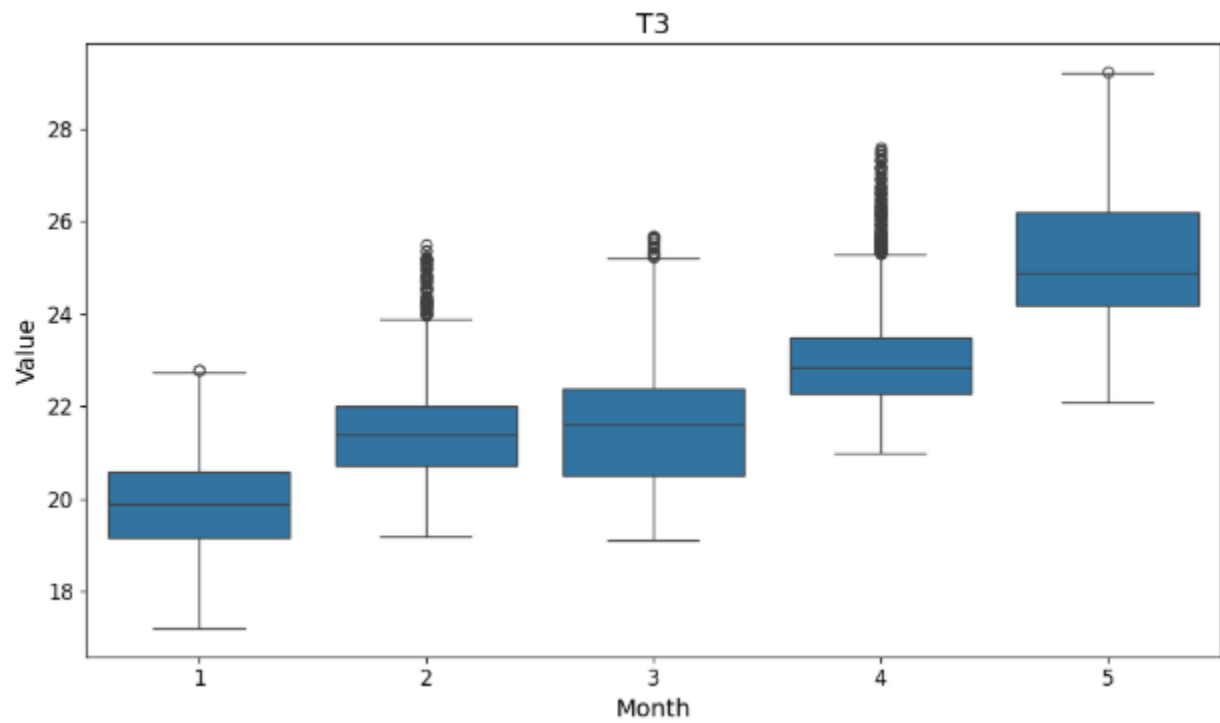


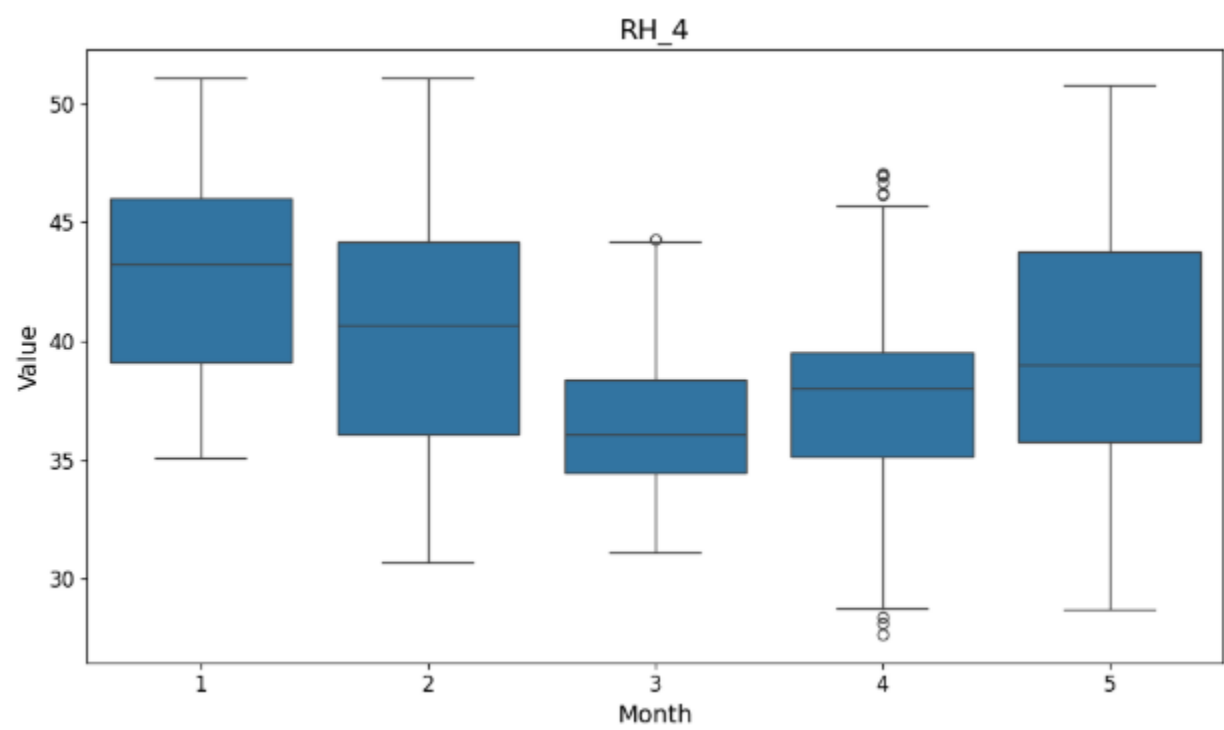
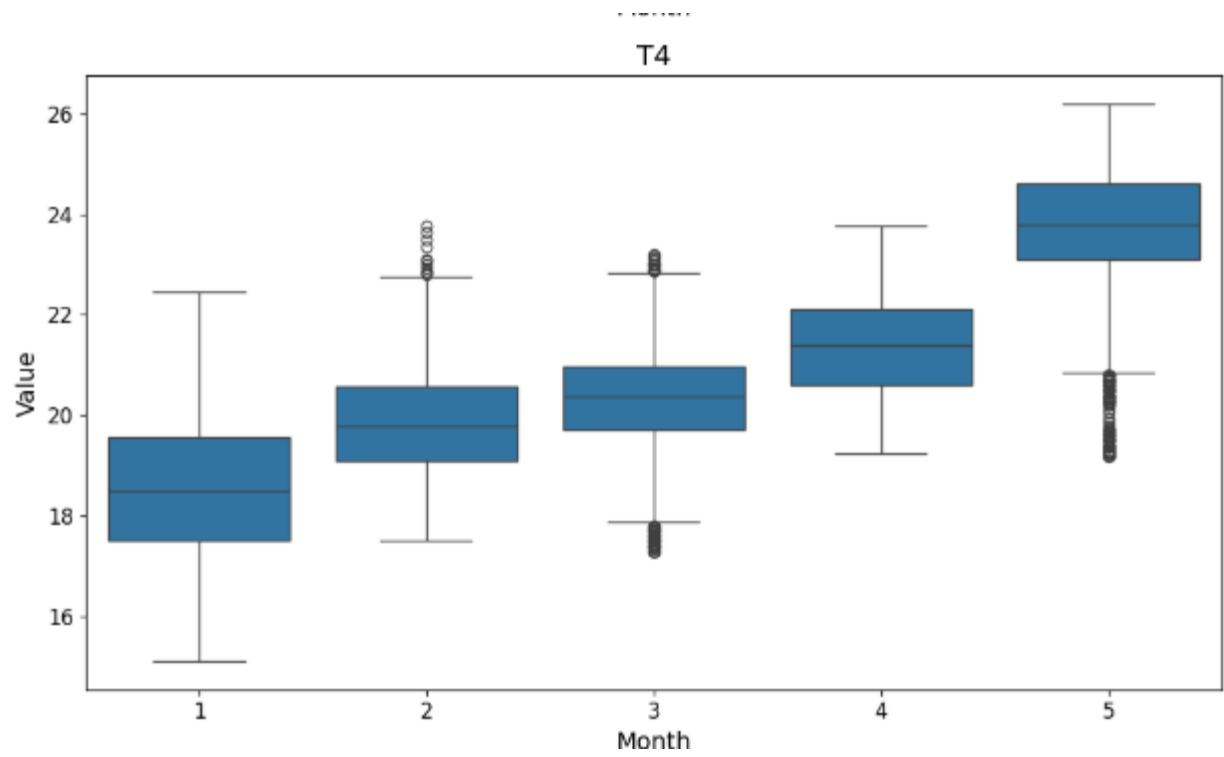


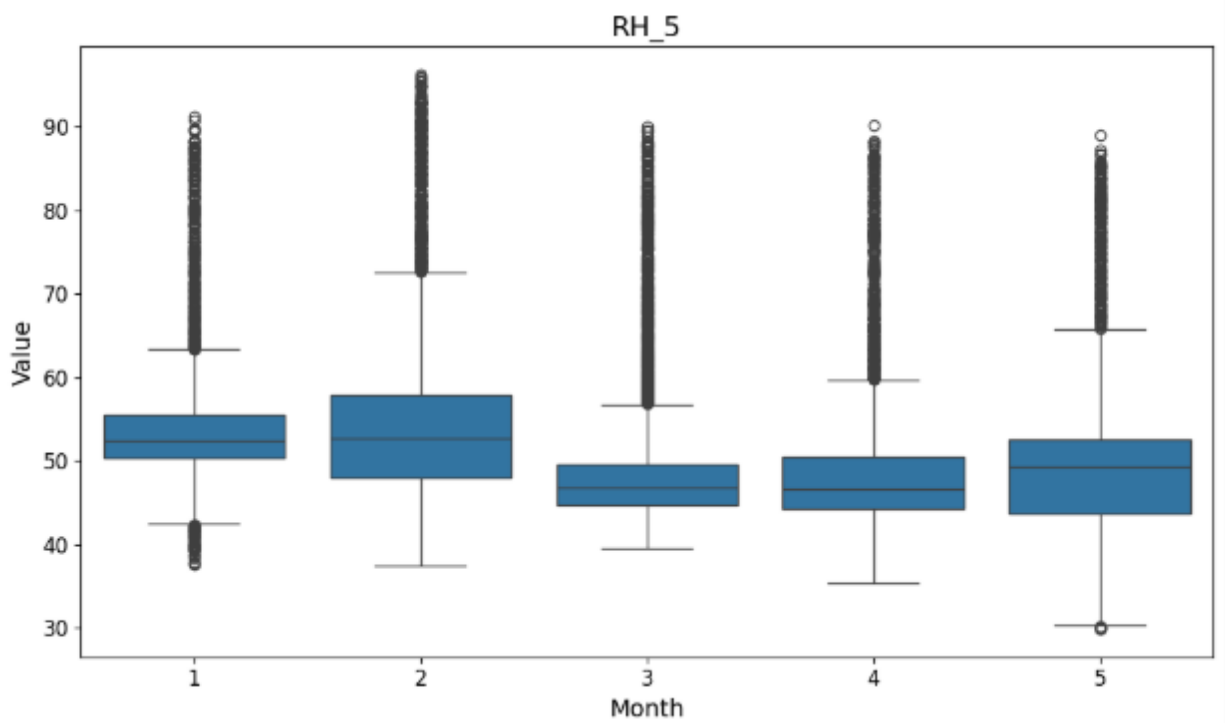
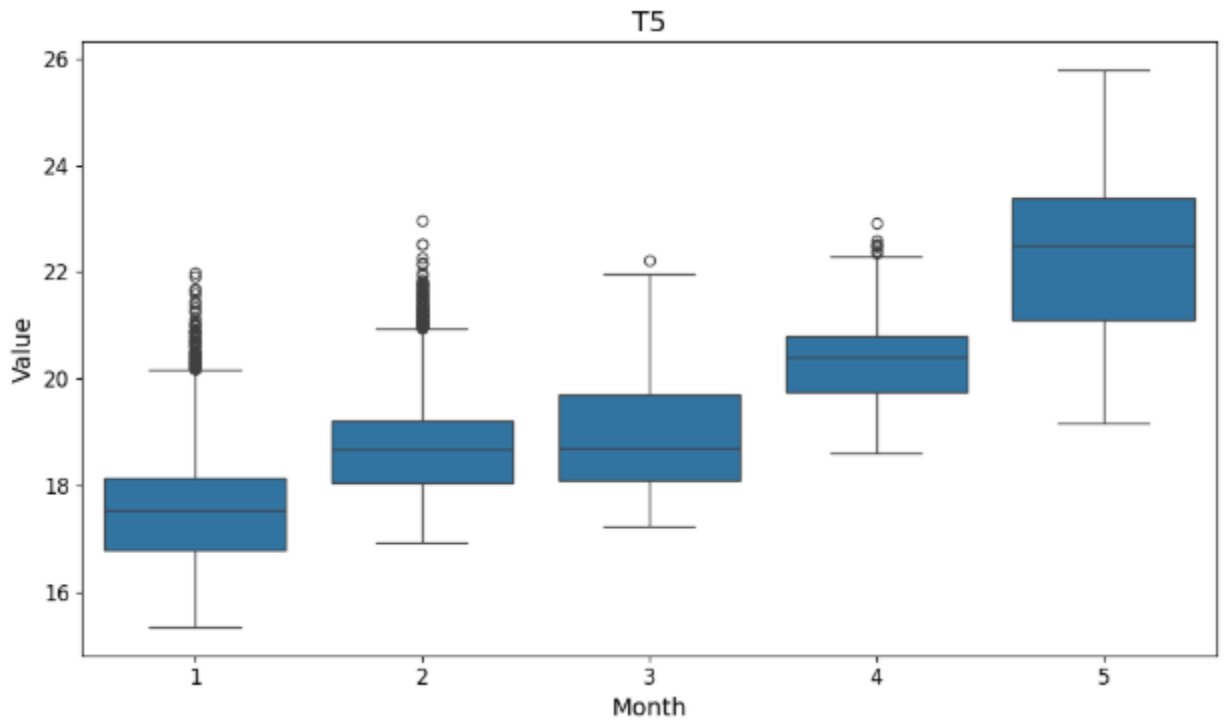


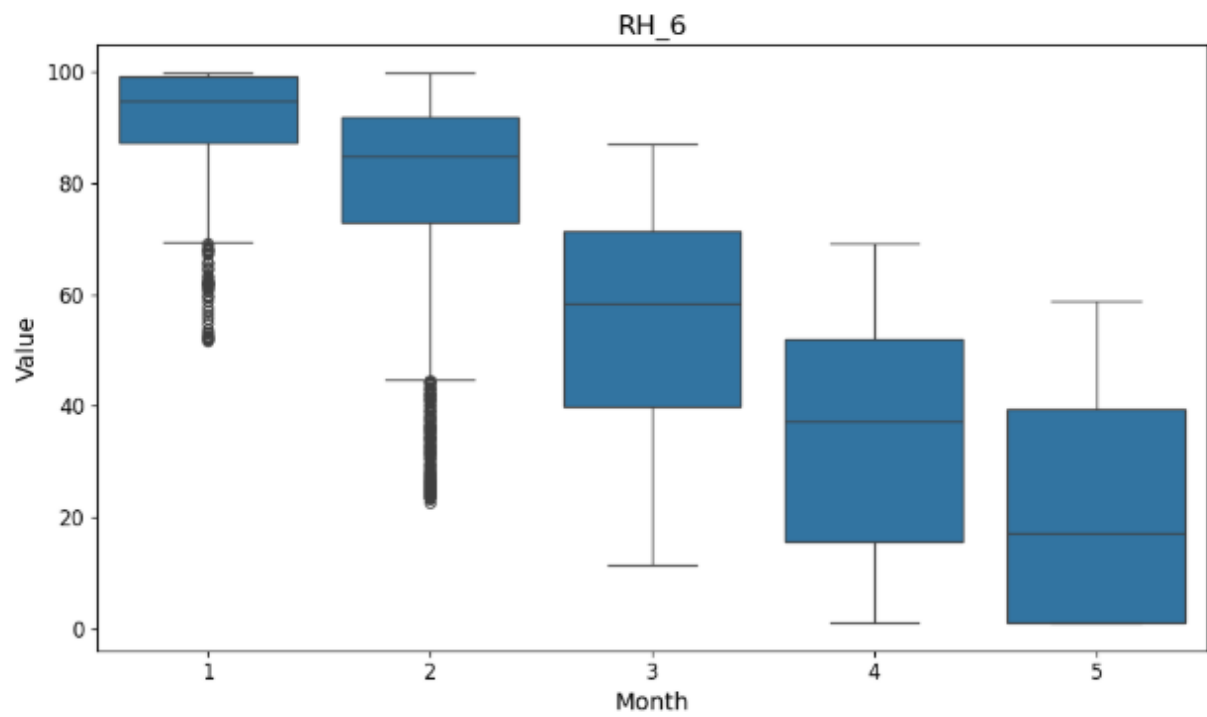
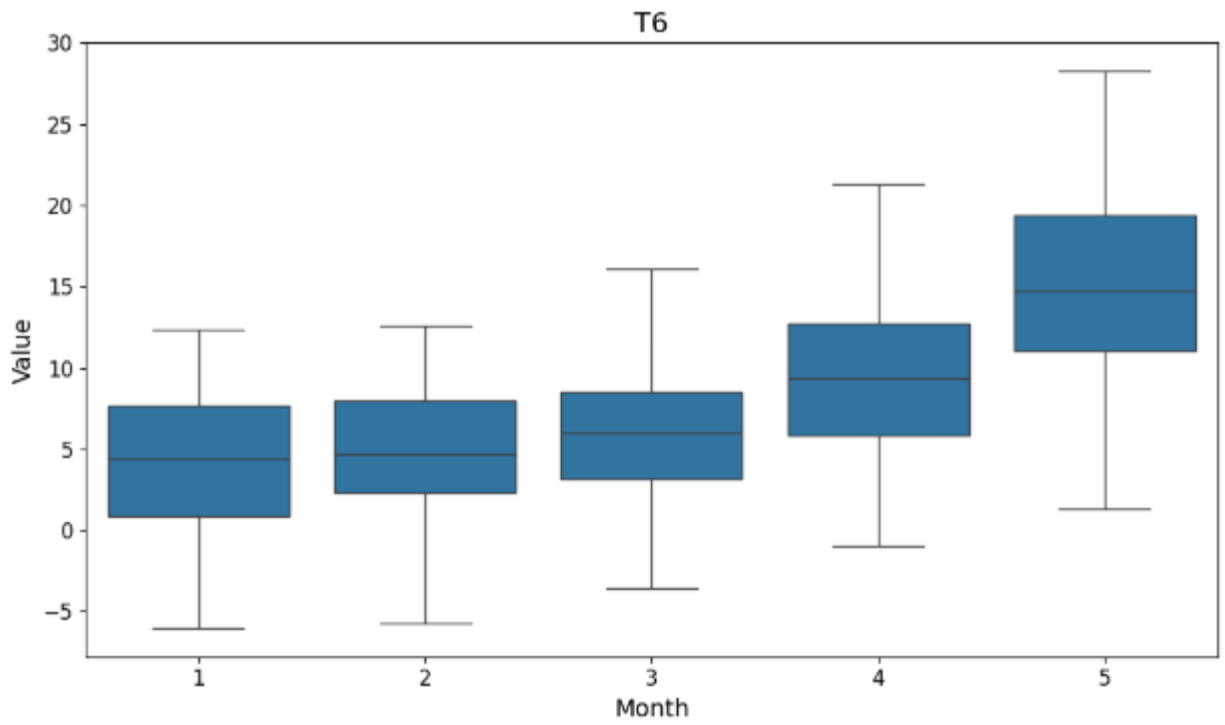


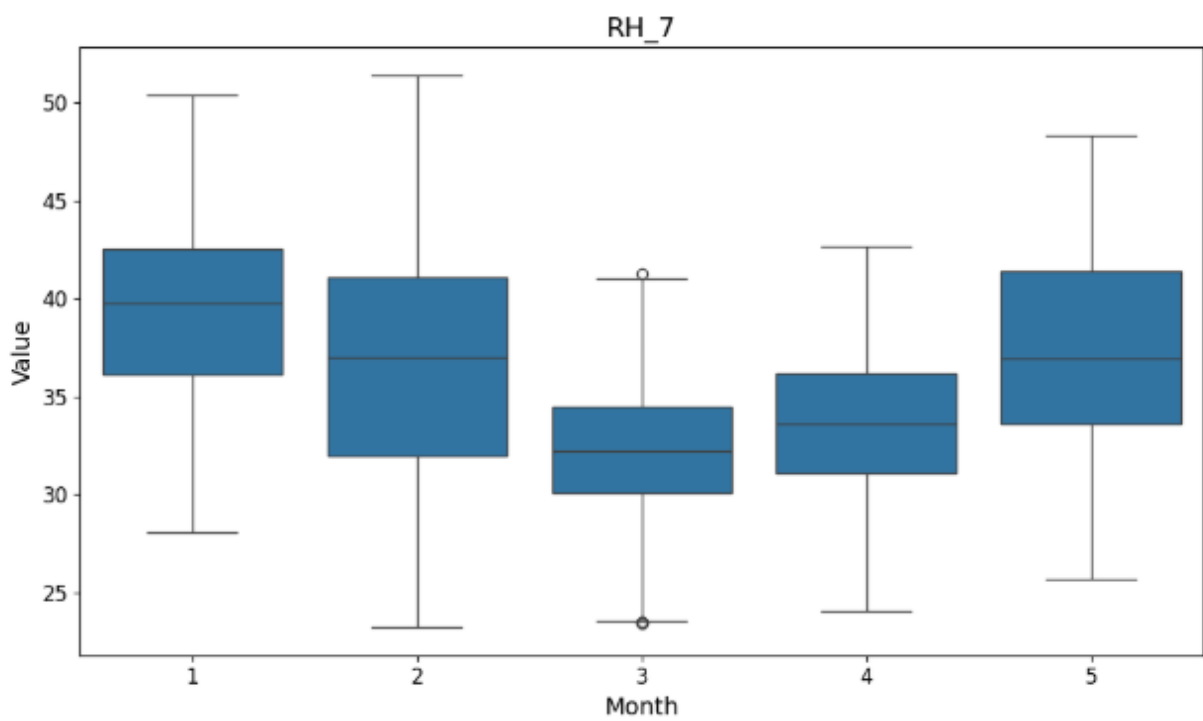
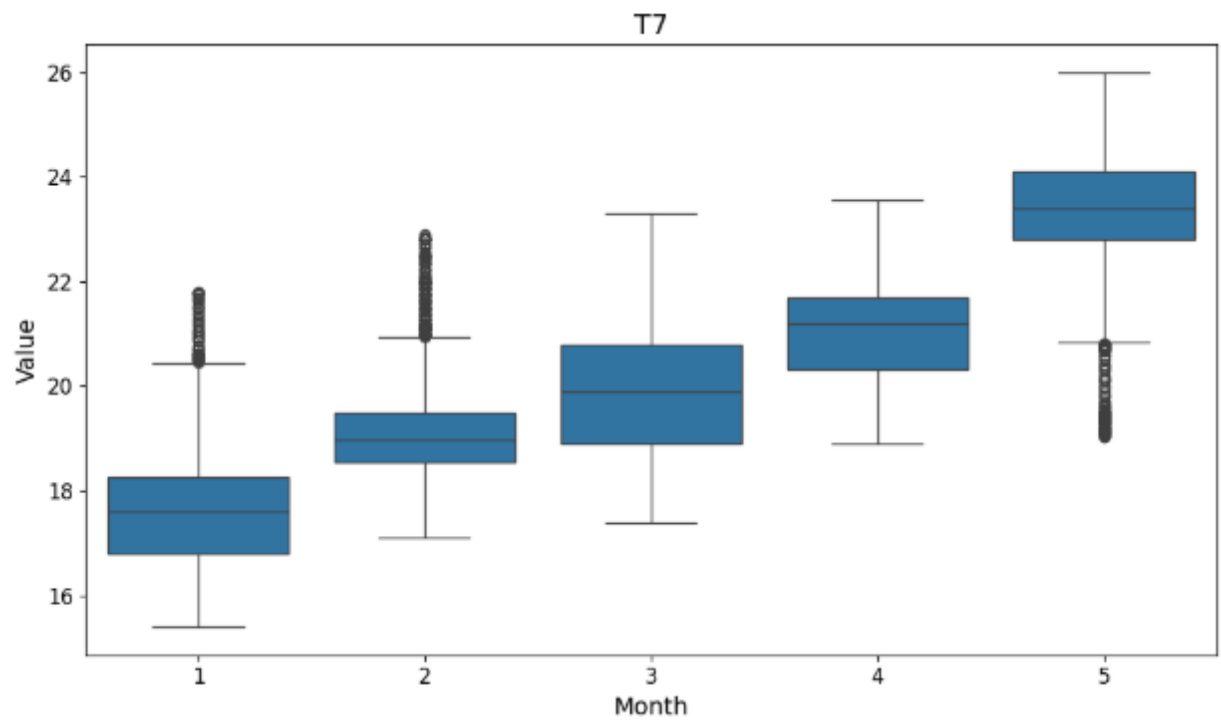


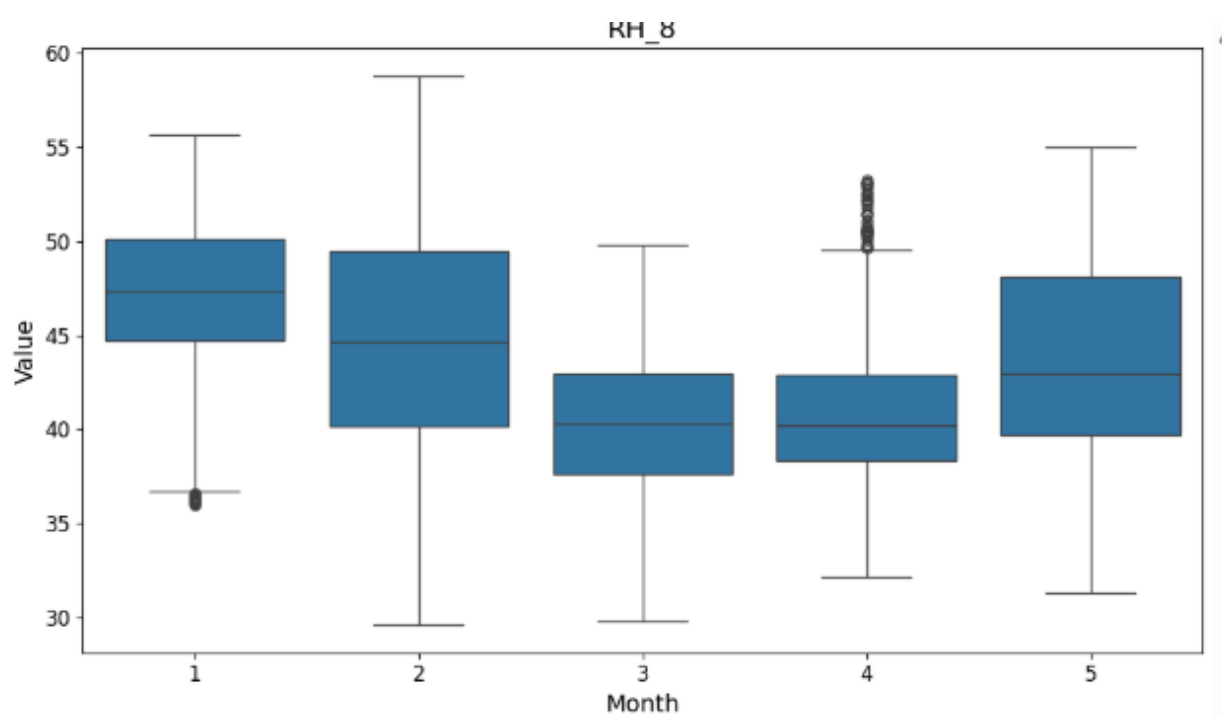
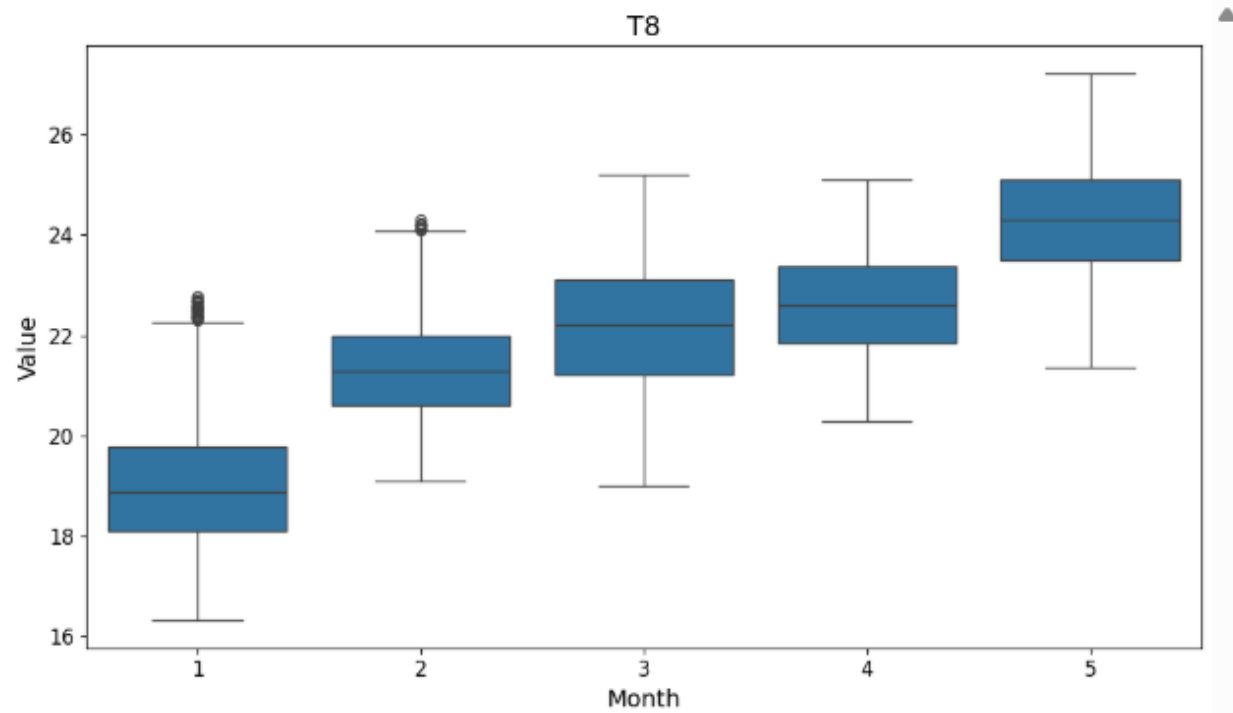


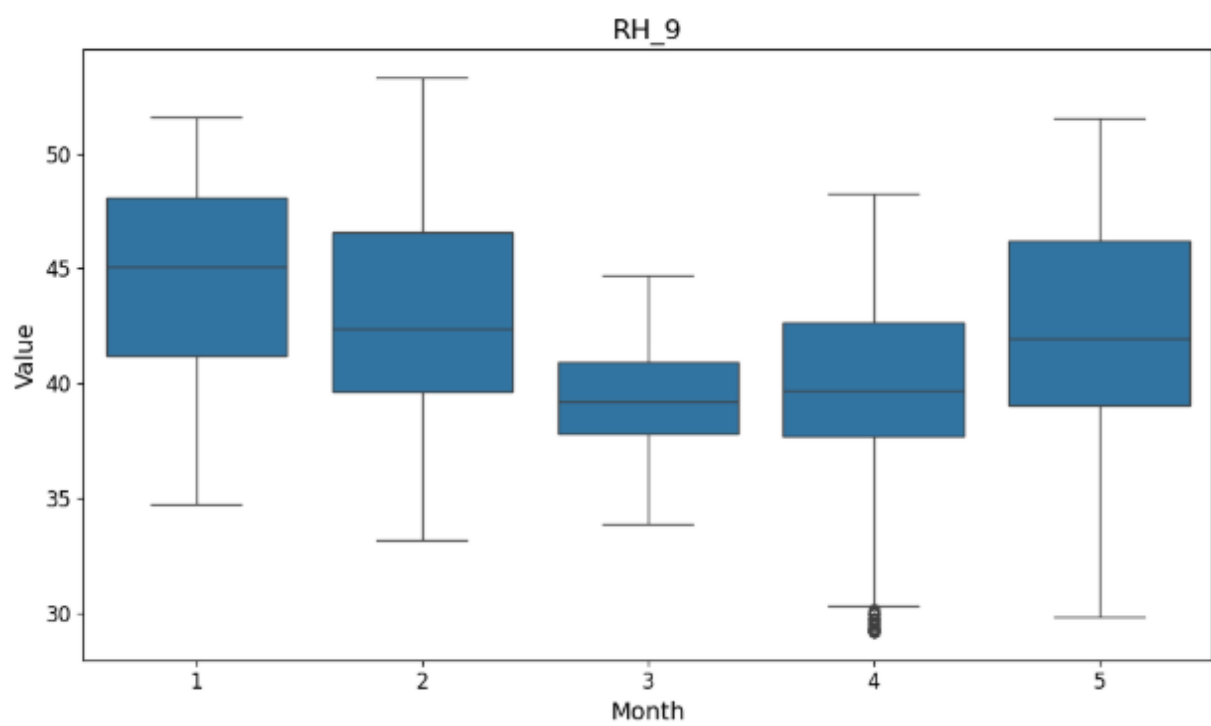
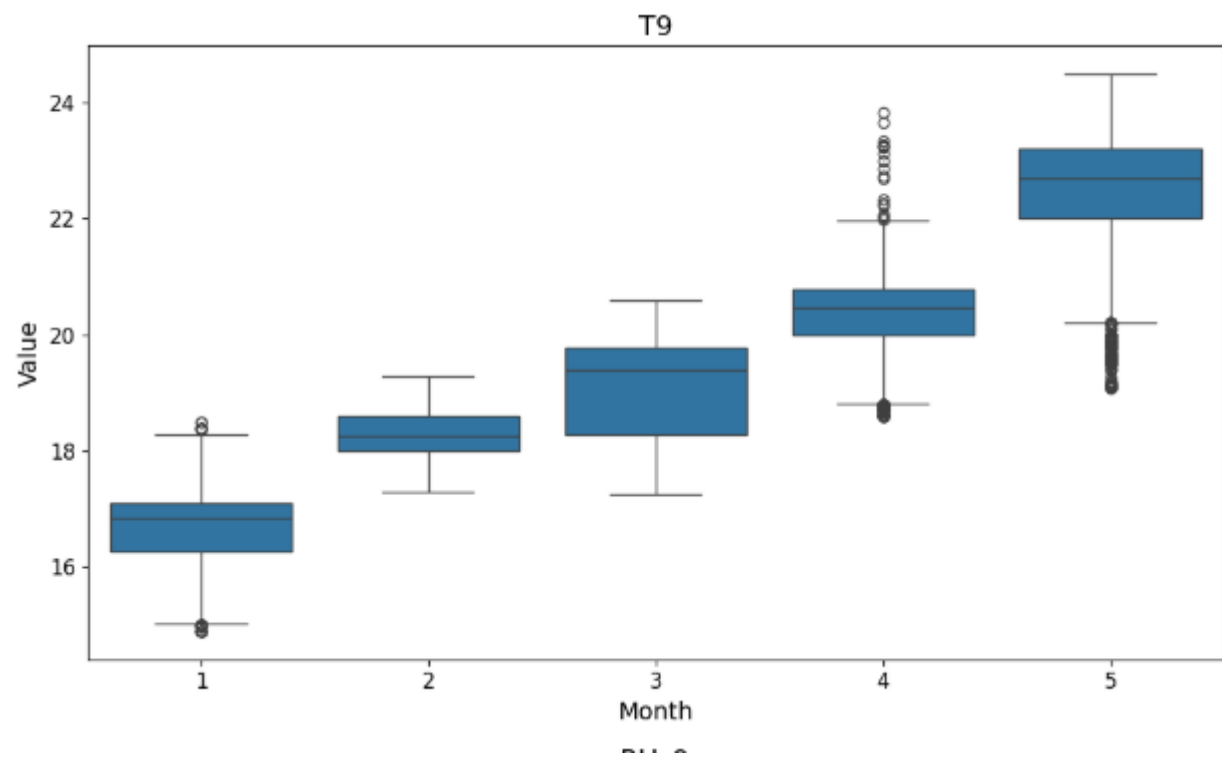


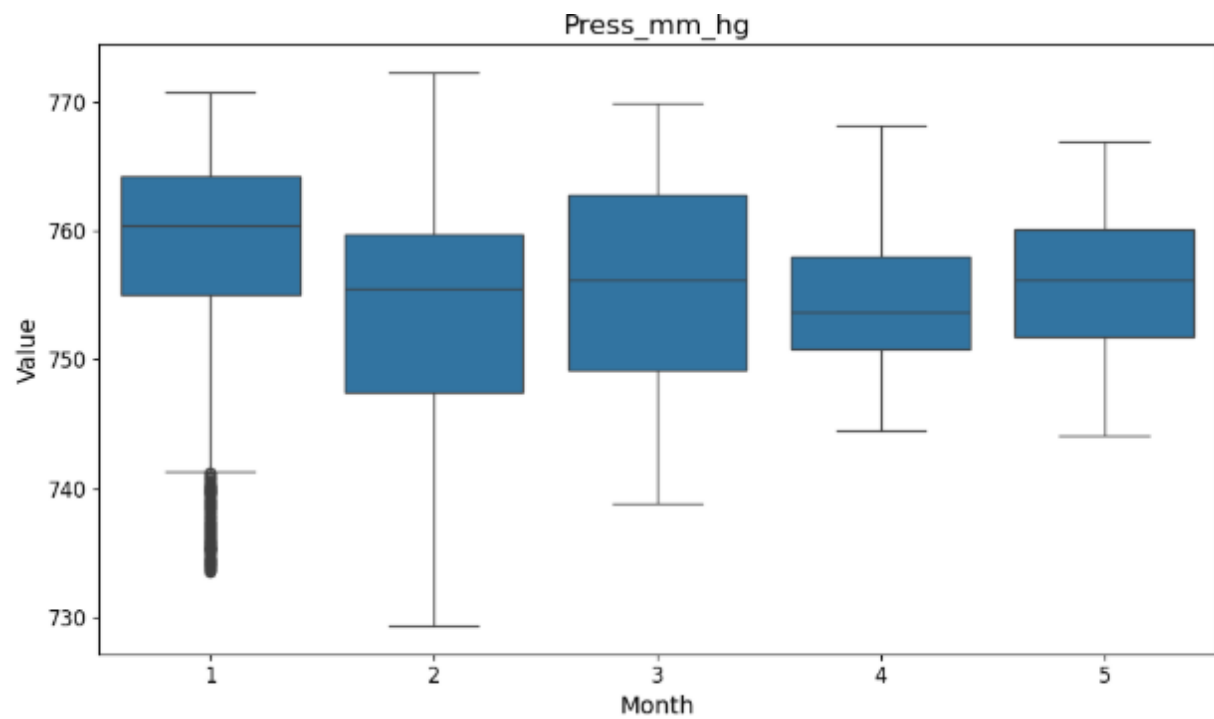
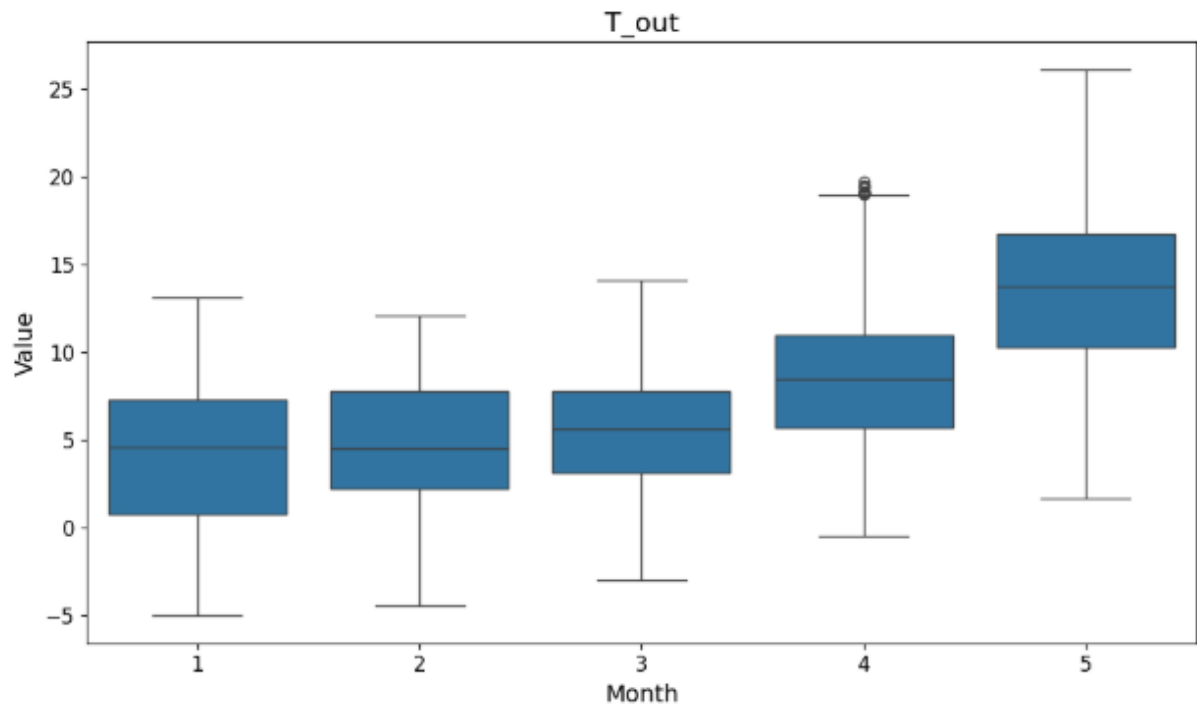


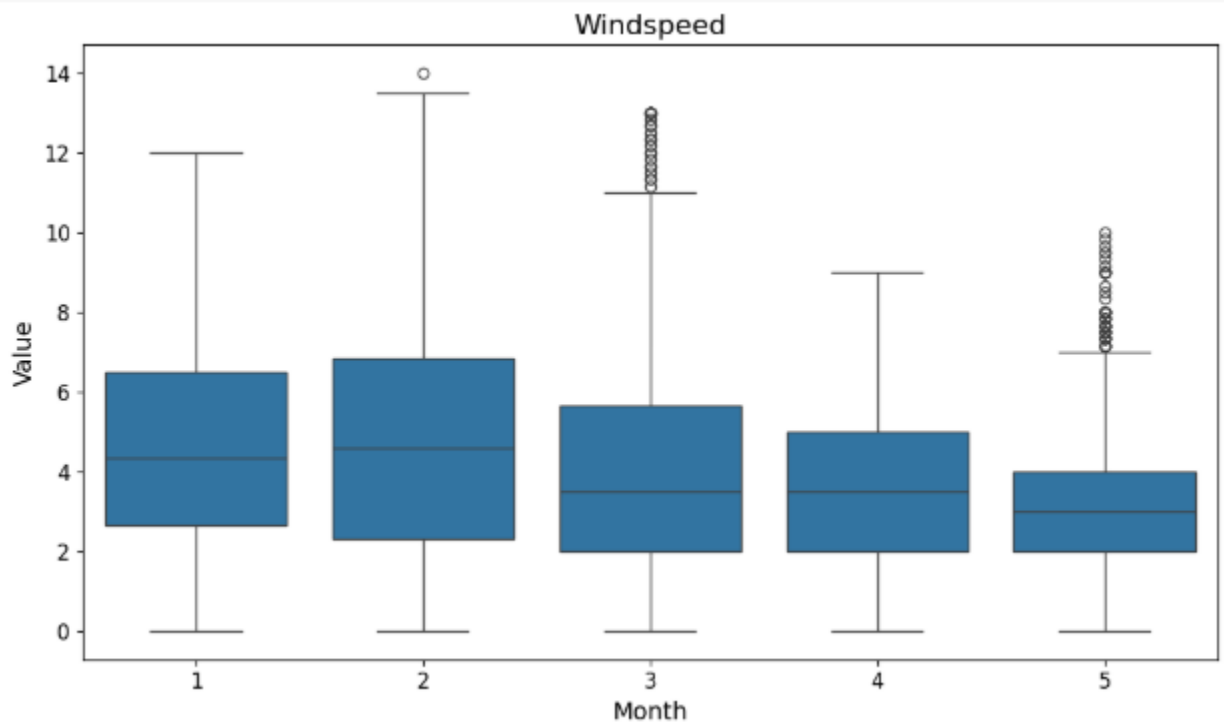
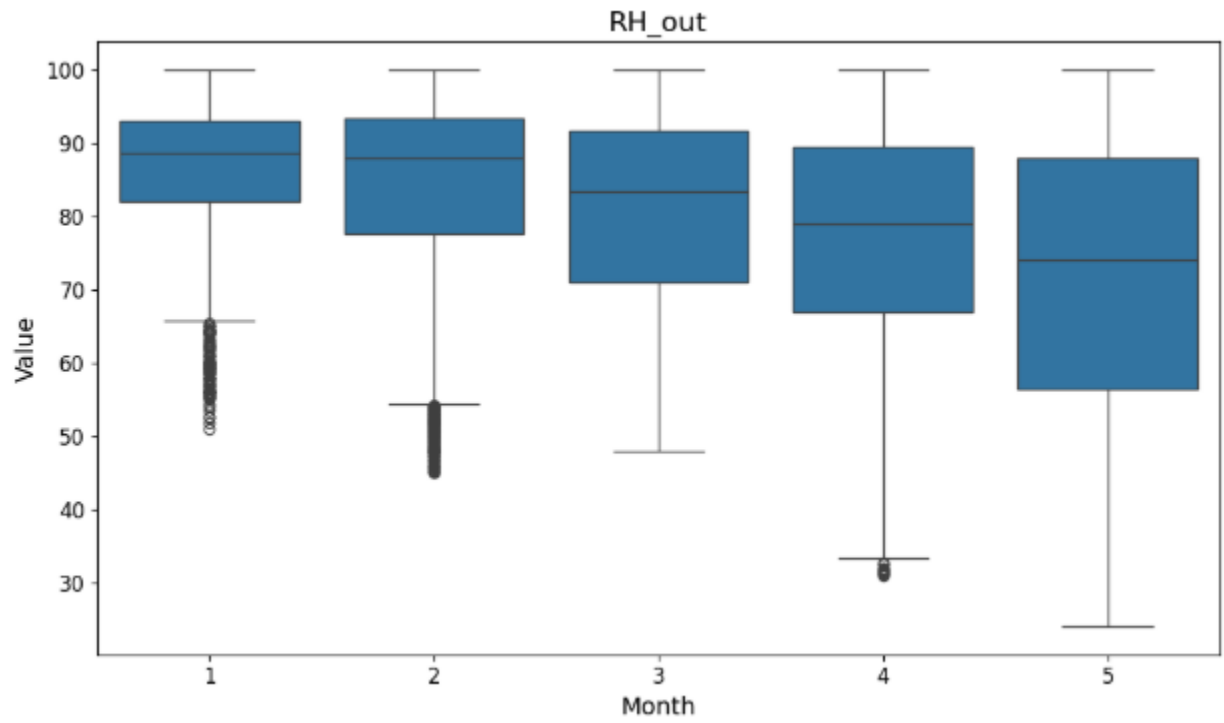


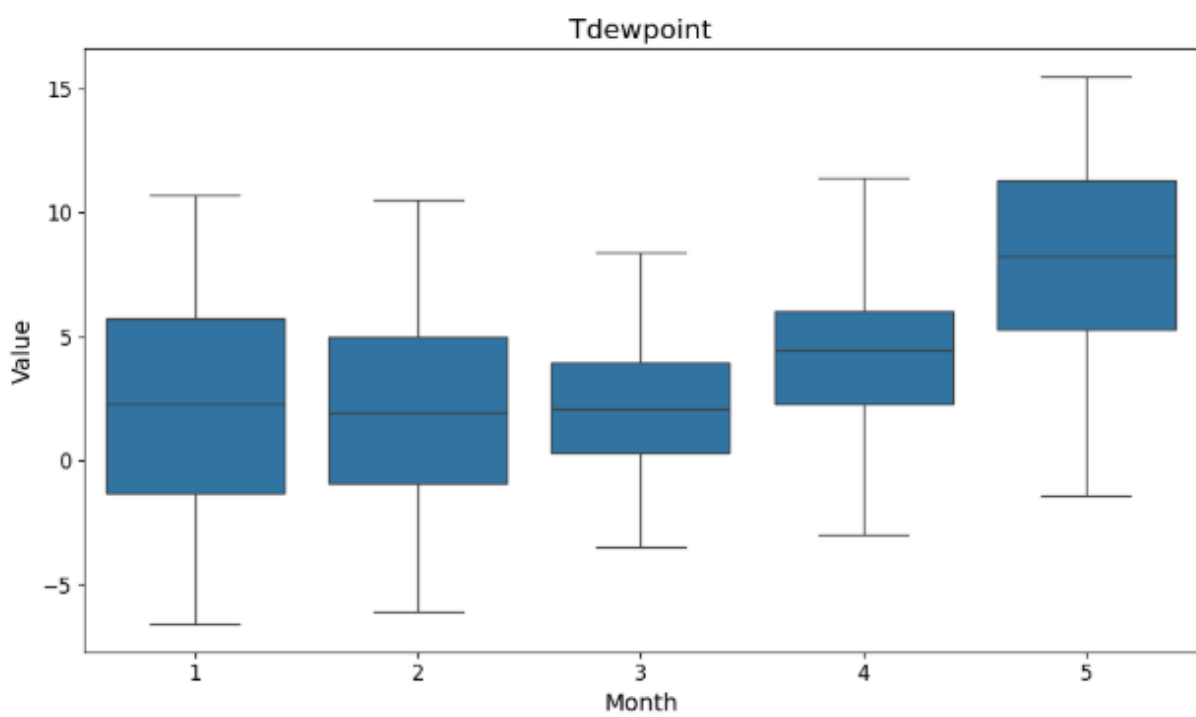
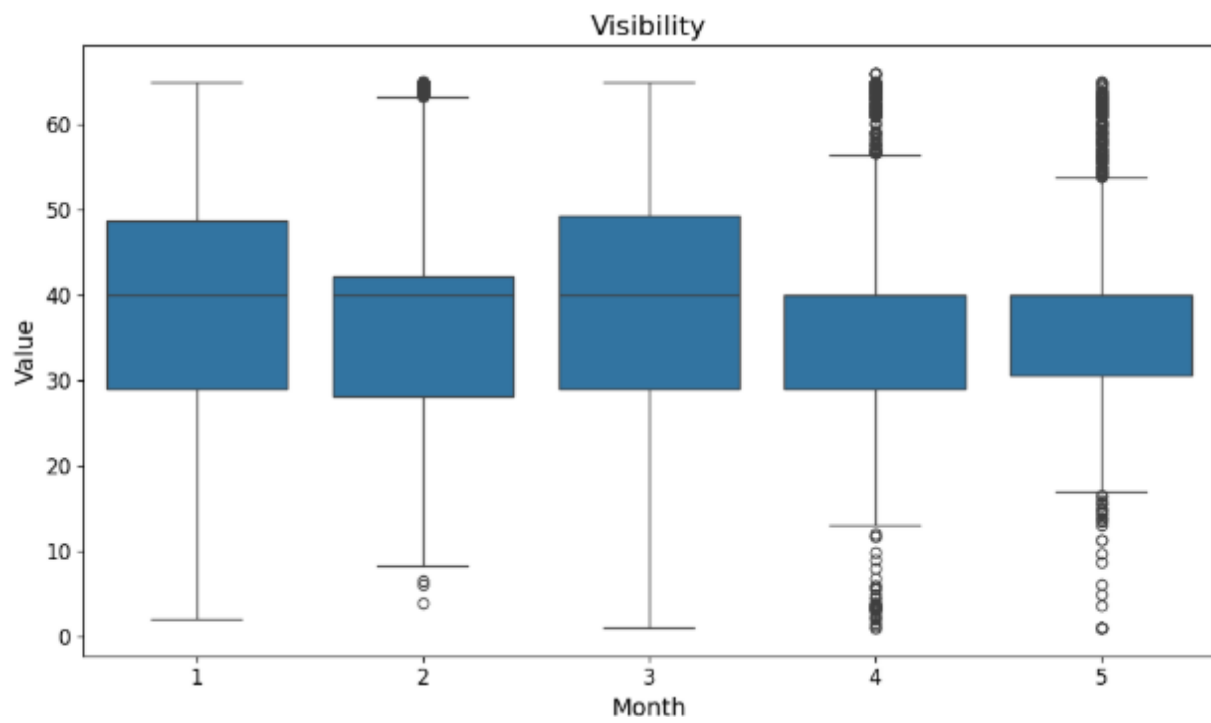


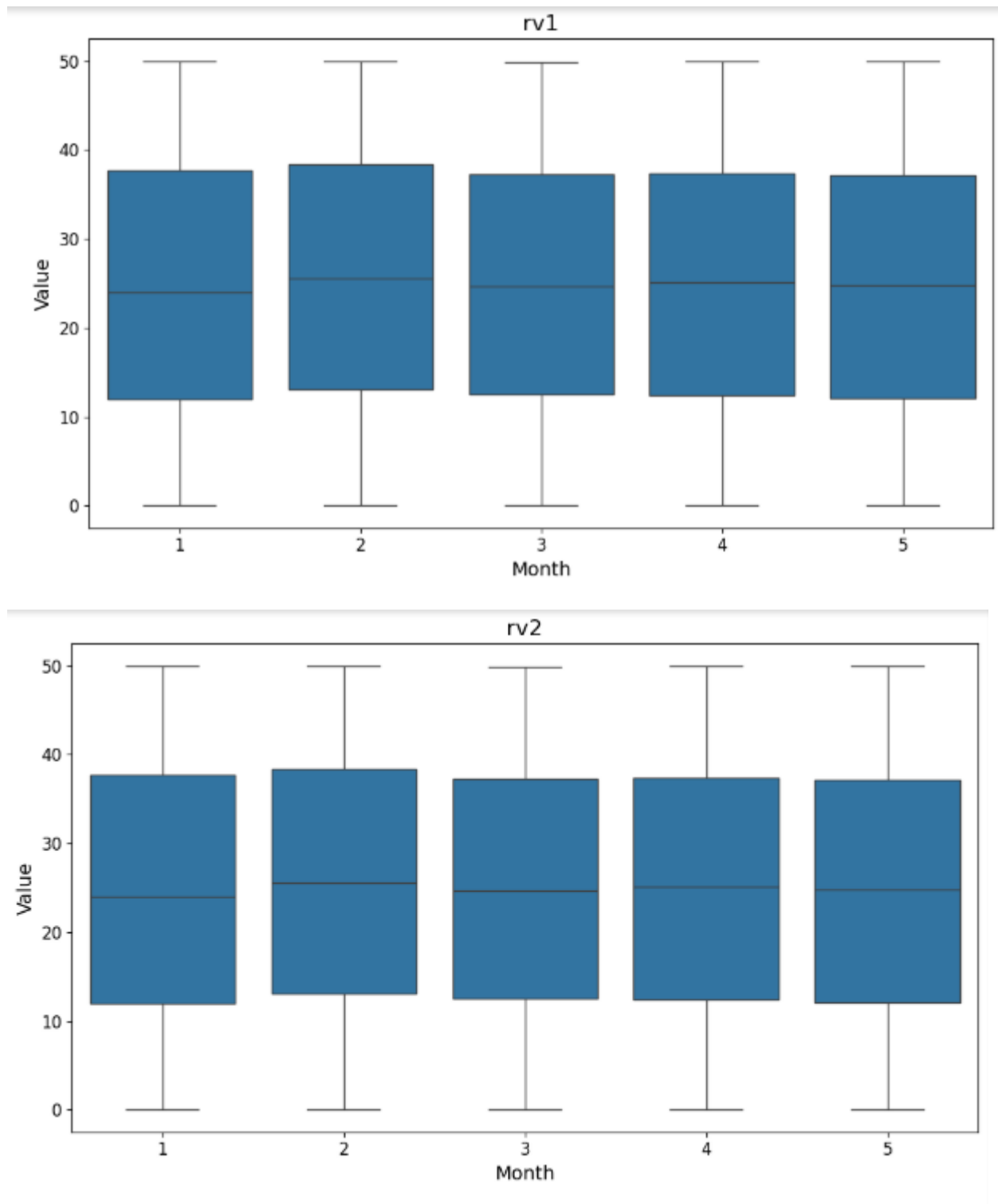






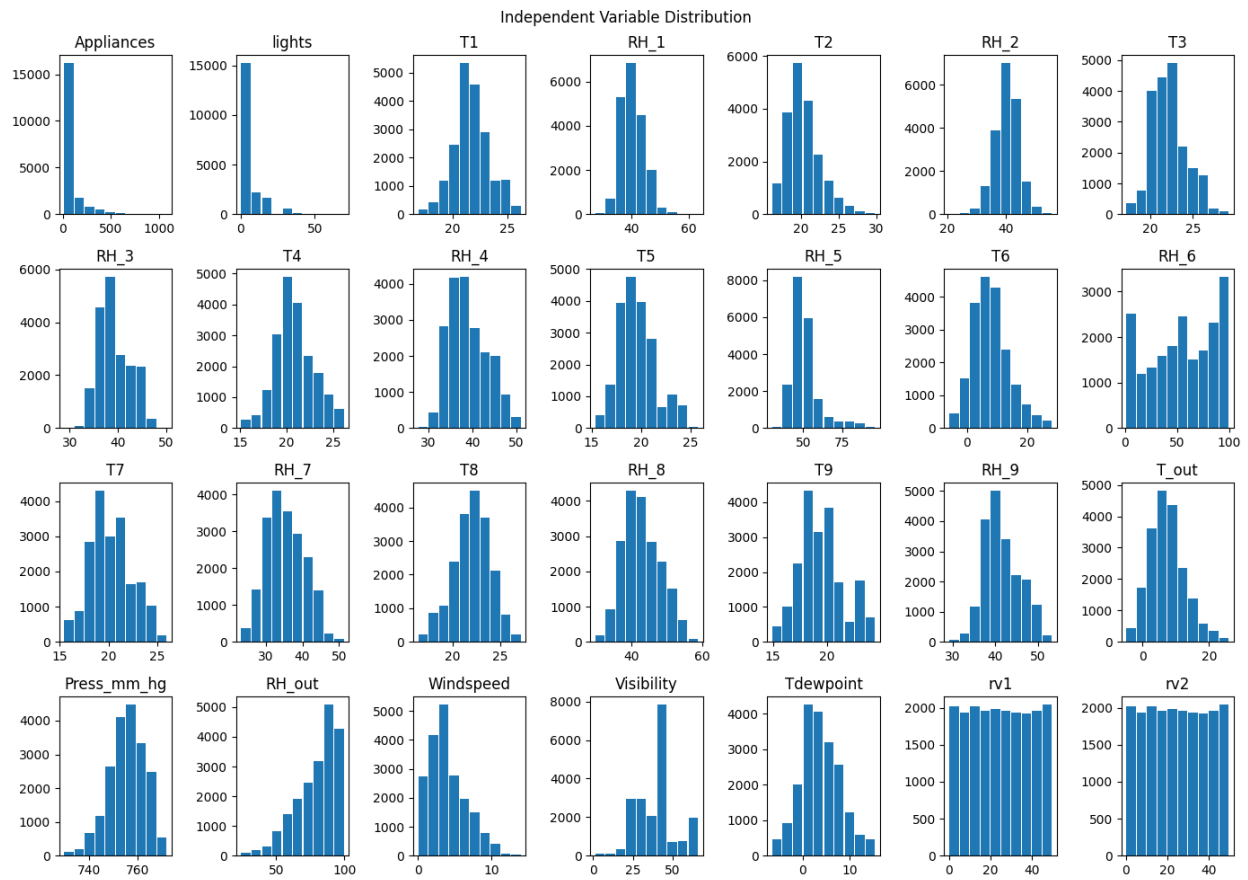






Observations:

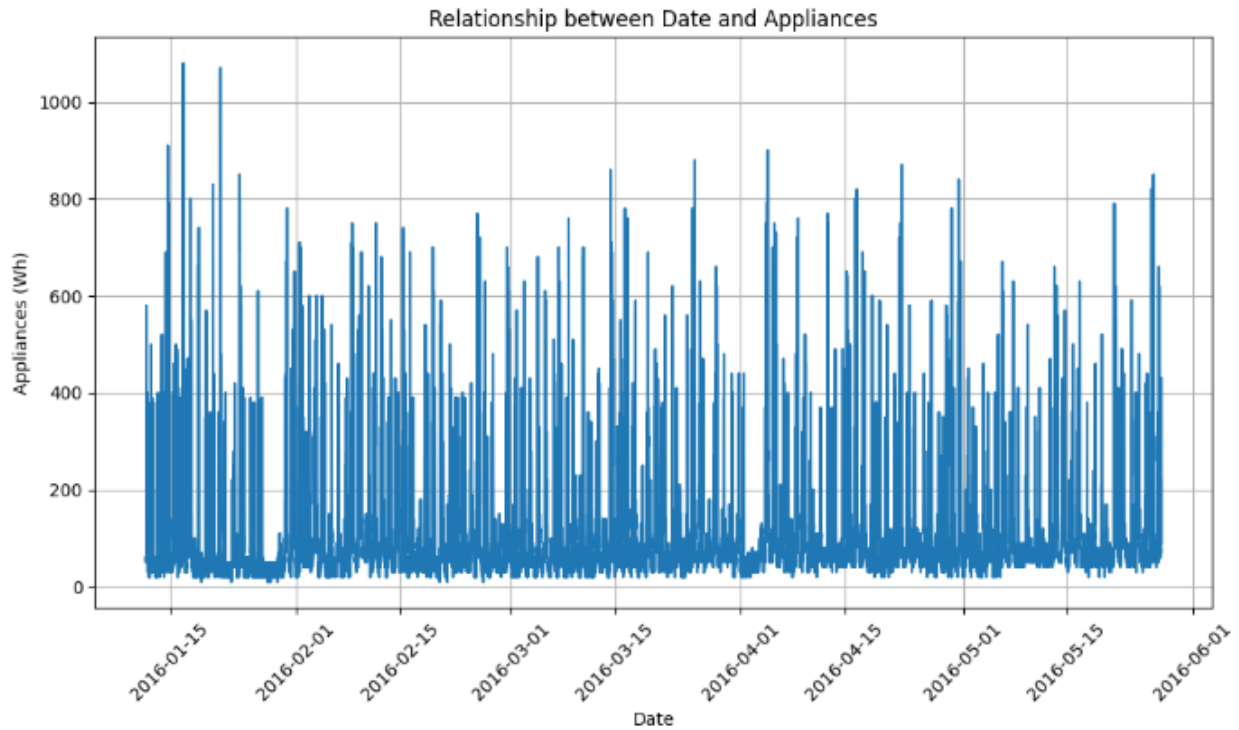
- For temperature (T), there is an overall increasing trend in the data.
- Other data show no significant features.



Observation:

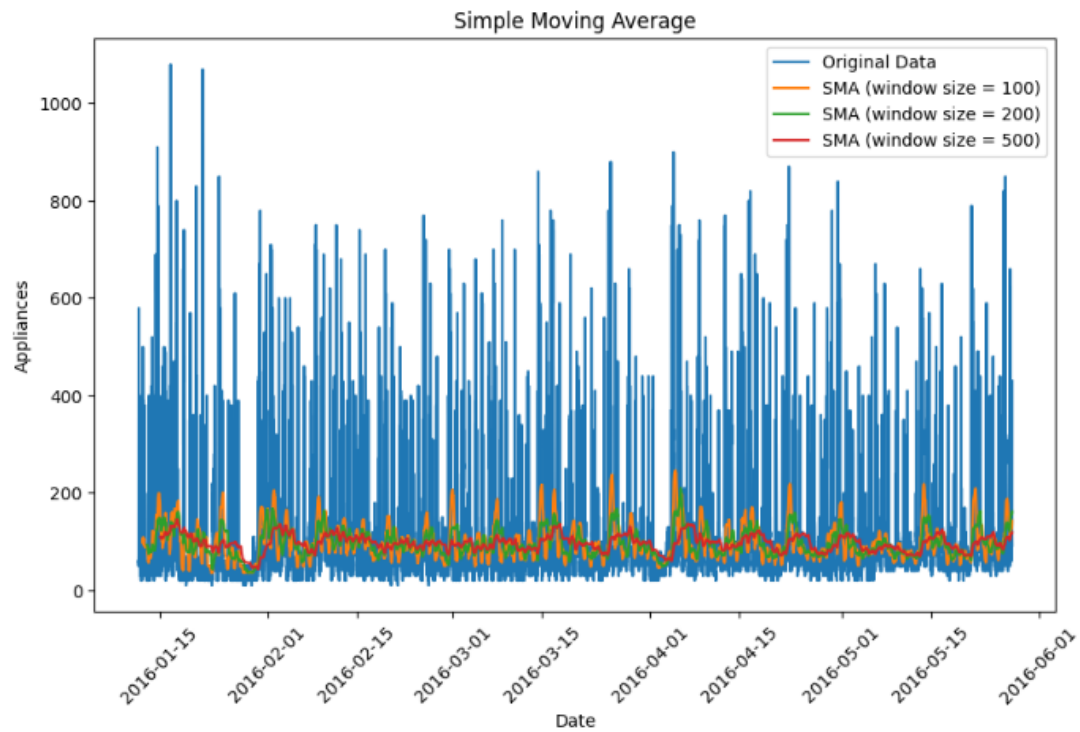
- Temperature and humidity exhibit a normal distribution, while other variables may be skewed or display multiple peaks.

A closer look at Appliances



Observation

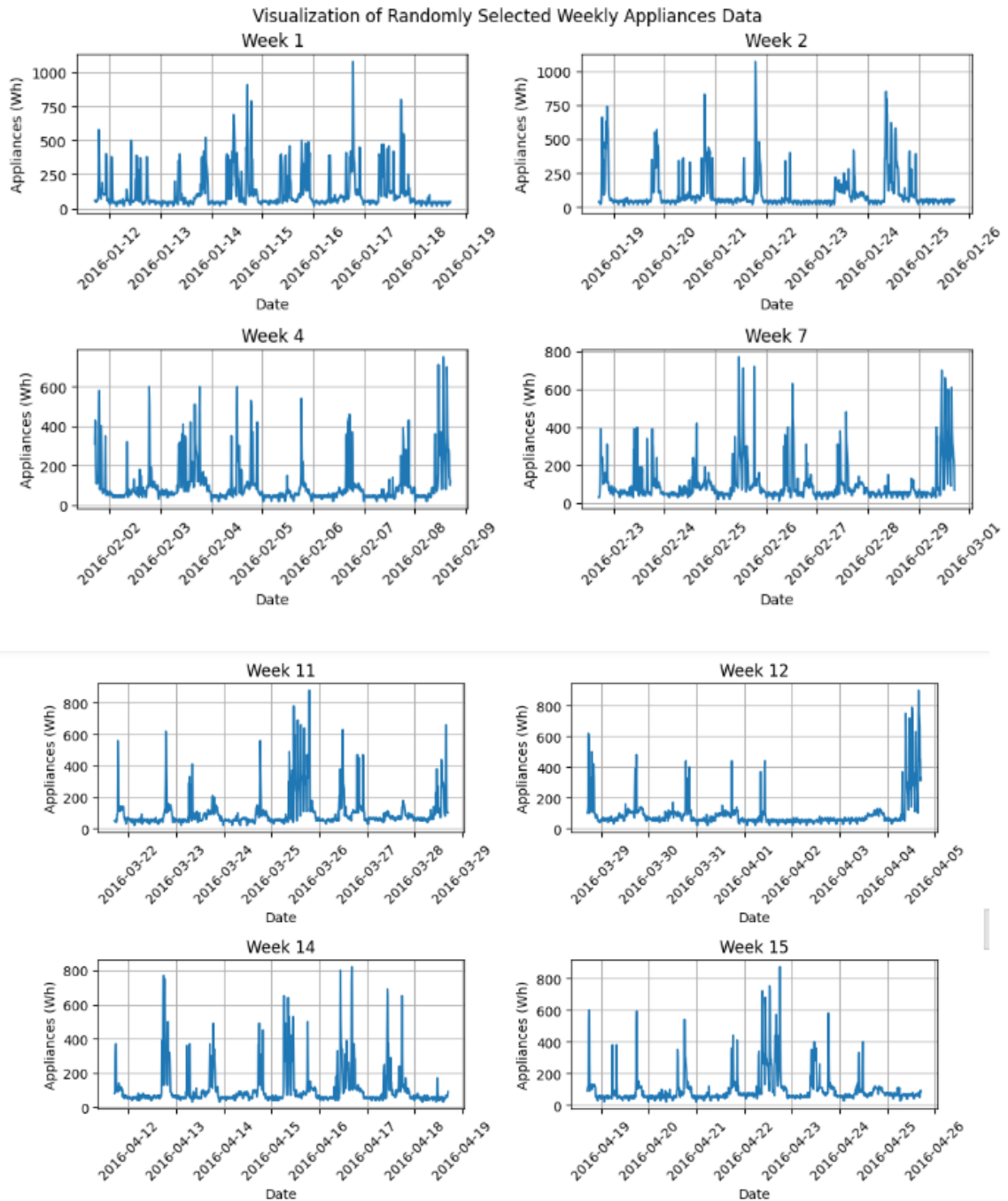
- From the plots, we can see that electricity usage (shown by the 'Appliances' variable) doesn't follow a clear or obvious pattern over time. But just because we can't easily see a pattern doesn't mean there isn't one. It means we need to look deeper and use data analysis techniques to find any hidden trends or relationships in the data.



Observation

- After using the Simple Moving Average (SMA) with window sizes of 100, 200, and 500 to smooth the data, the graphs didn't show any clear trends. This means that the 'Appliances' usage stayed mostly steady over time, without big changes or patterns.

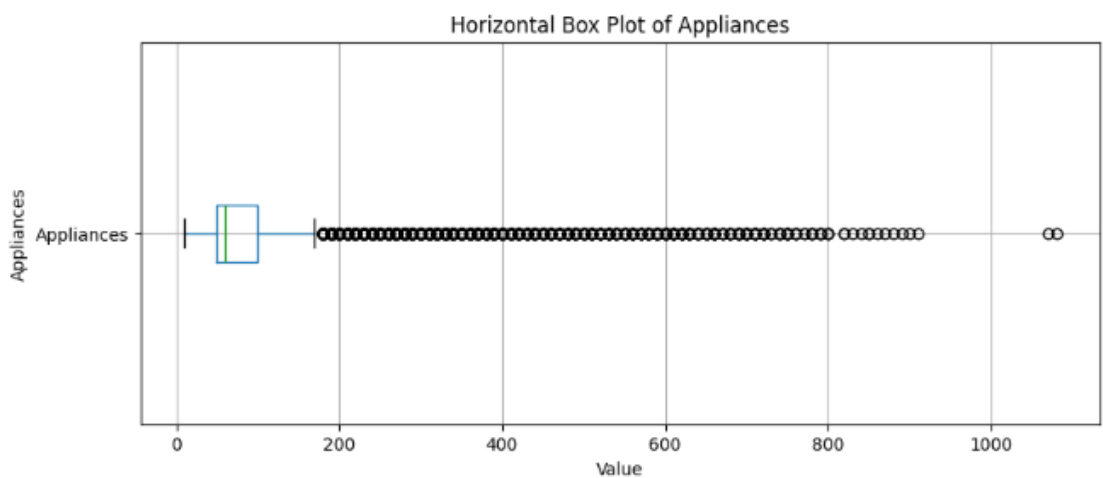
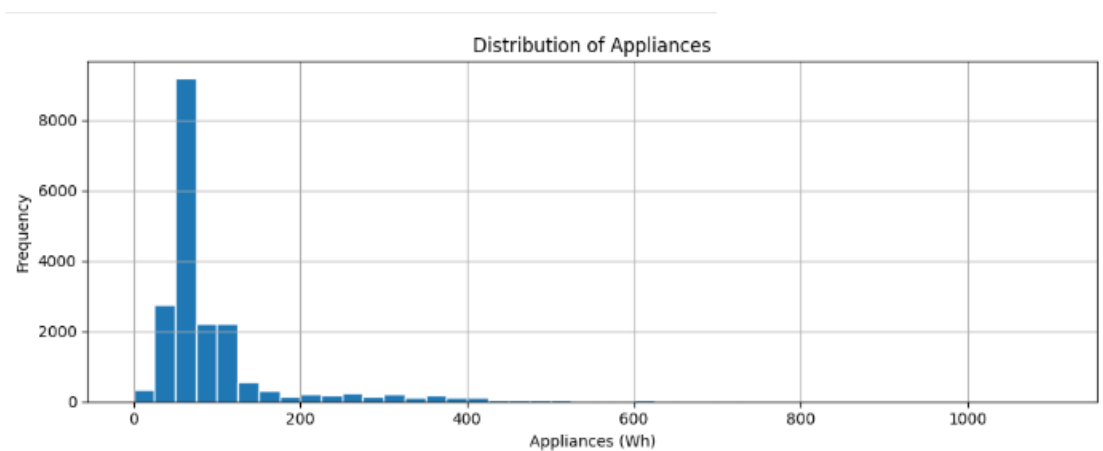
Visualization of Randomly Selected Weekly Data



Observation

- When looking at a random week of data, we noticed a clear pattern showing regular ups and downs in energy use. Even though the week was picked at random, similar trends kept appearing, which means there are hidden habits or routines in how appliances are used. This matches what we saw earlier in the hourly analysis.

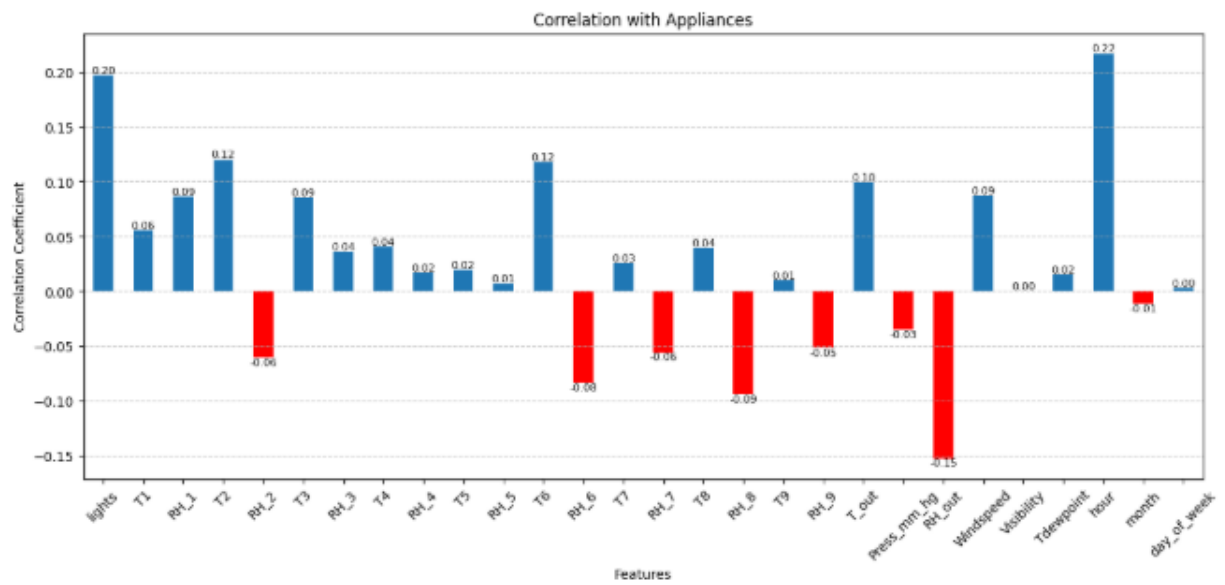
Plotting the histogram for the "Appliances" column



Observation

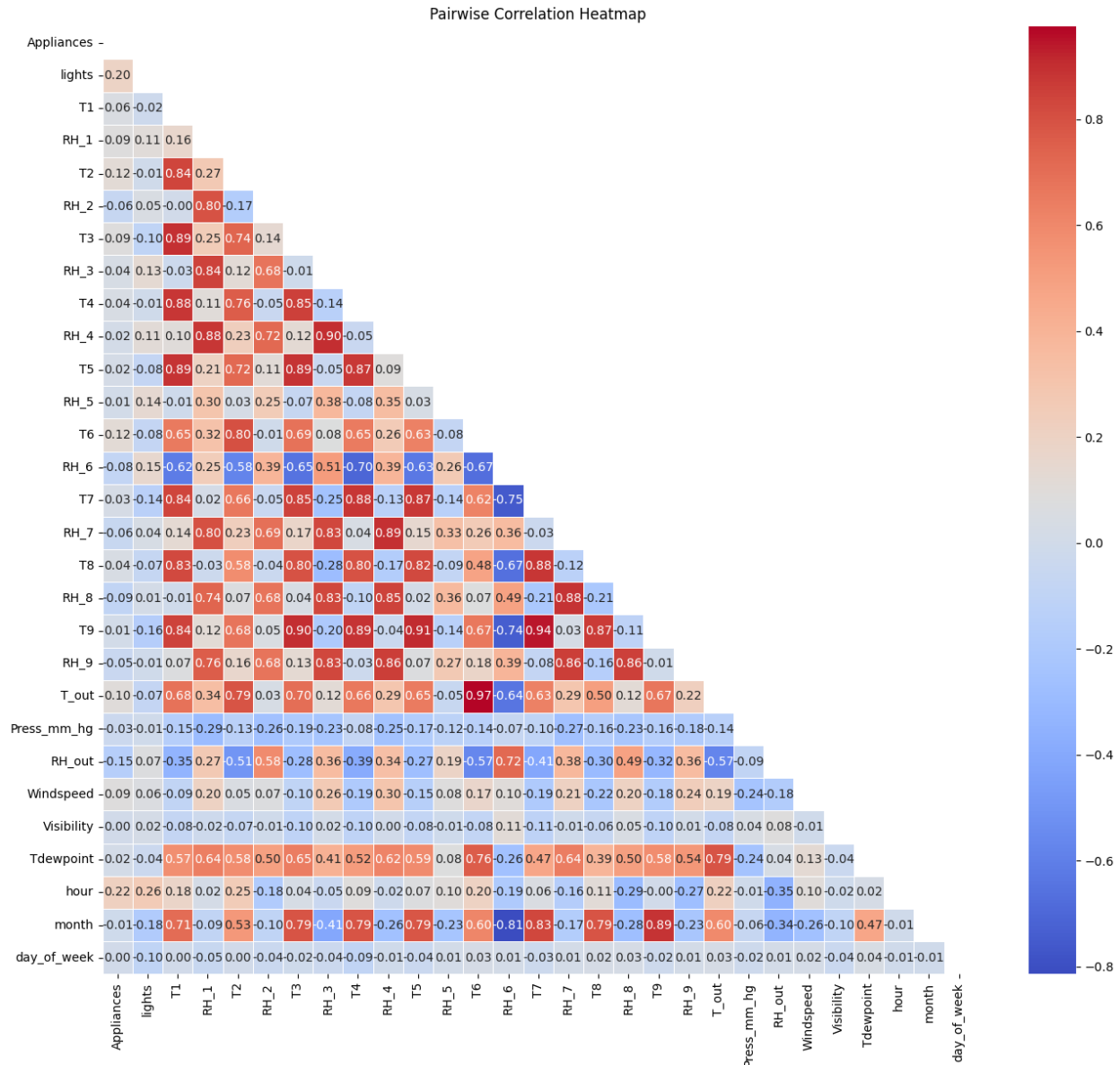
- The box plot has a long right-hand side tail, it indicates that it is positively skewed with some high energy use.

Correlation Analysis



Observation

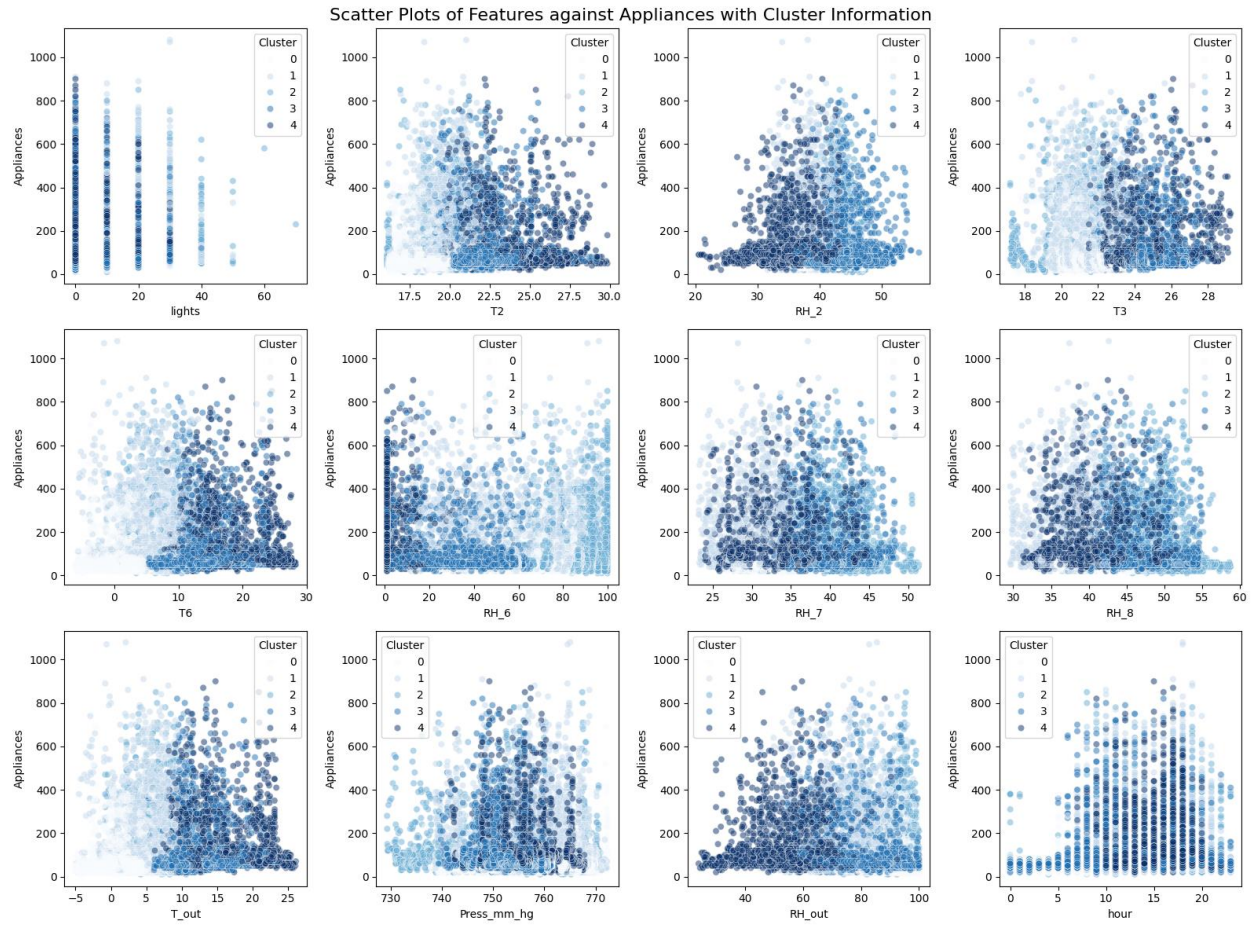
- The time of day shows the strongest link to appliance electricity use (correlation: 0.22), followed by the number of lights on (0.20). Air pressure has a weak negative correlation (-0.15), suggesting that appliance usage slightly decreases as air pressure rises.



Observation:

- This heatmap shows how different features are related to each other. We can see that temperatures are strongly linked with each other, and humidity levels also show some positive relationships, but not as strong.
- One interesting point is RH_6, which measures outdoor humidity. It has a strong negative relationship with temperature — meaning that when the temperature goes up, outdoor humidity tends to go down.
- On the other hand, our target feature, 'Appliances', doesn't show a strong relationship with any of the temperature or humidity features. So, these factors might not directly affect how much electricity appliances use.

Clustering Analysis



Observation

- The clusters in the data are clearly separated, especially in temperature-related plots, which show distinct groupings. Humidity-related plots show less clear separation. There's also a visible pattern related to lights. Overall, the data seems to be grouped well, making the classification look reliable.

Preprocessing

```
[ ] #Checking for Duplicate Rows
duplicate_count = dataset[dataset.duplicated()].shape
duplicate_count
```

➡ (0, 29)

- All rows are unique, and there are no duplicate rows in the dataset.

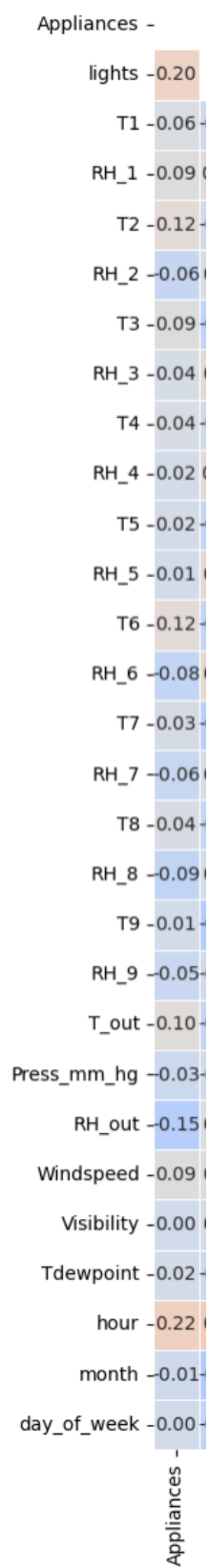
```
#Missing Values/Null values Count
dataset.isnull().sum()
```

- There are no missing values in any of the columns.

```
# 4. Normalize the dataset using MinMaxScaler:
scaler = MinMaxScaler(feature_range=(0, 1))
scaled_lstm_data = scaler.fit_transform(lstm_data_values)
```

Feature Engineering

The variables `Visibility` (-0.00) and `day_of_week` (-0.00) have near-zero correlation with the target (`Appliances`), indicating no meaningful linear relationship. Features with very low correlation add little predictive value and are typically excluded to simplify the model and avoid noise. RV1 and RV2 (likely these variables) were not selected for this reason.



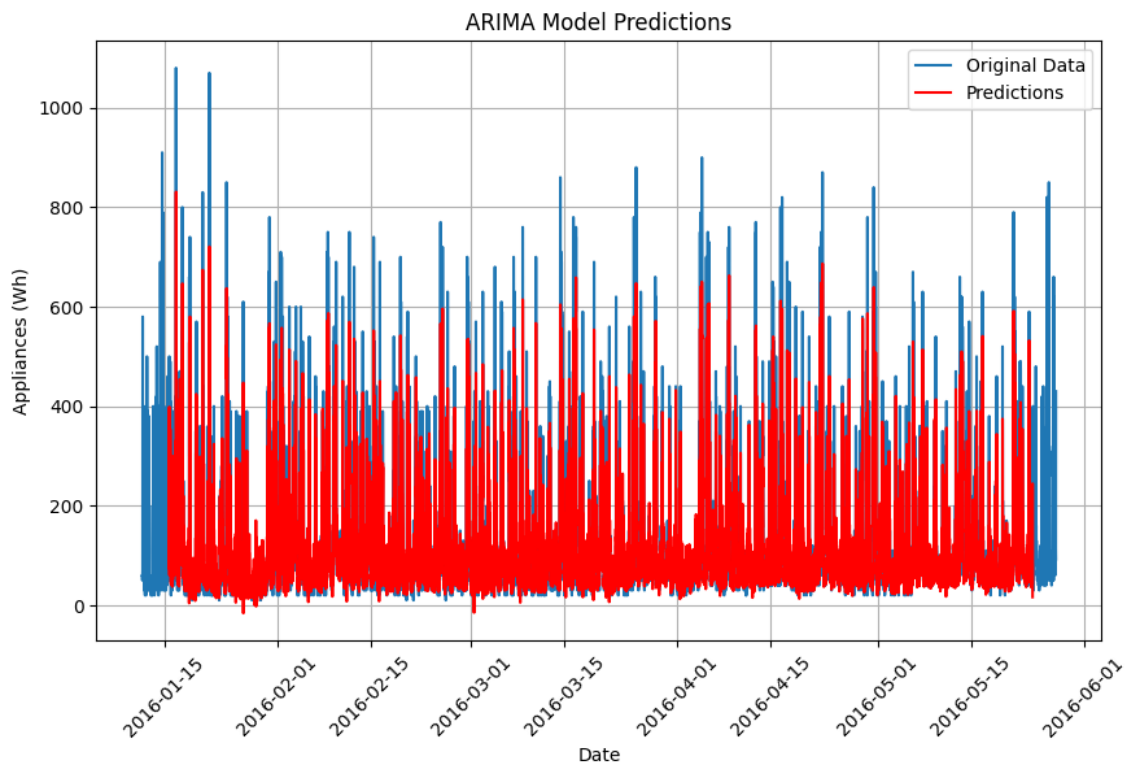
Model Design

The deep learning model used was an **LSTM-based sequential model** designed to capture temporal dependencies in energy usage data. The architecture included one LSTM layer with 64 units, followed by a Dropout layer to prevent overfitting, and a Dense output layer.

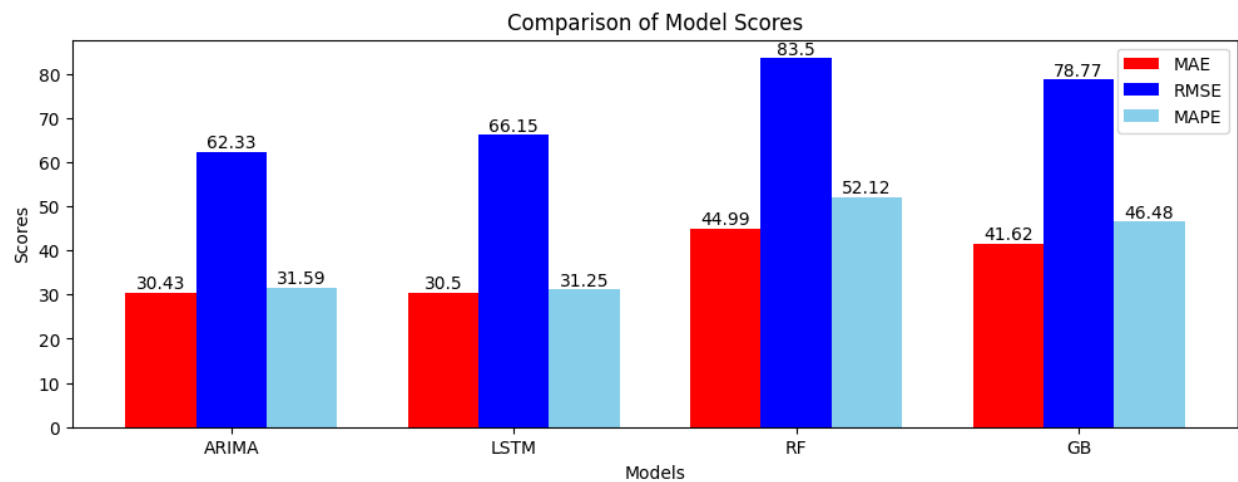
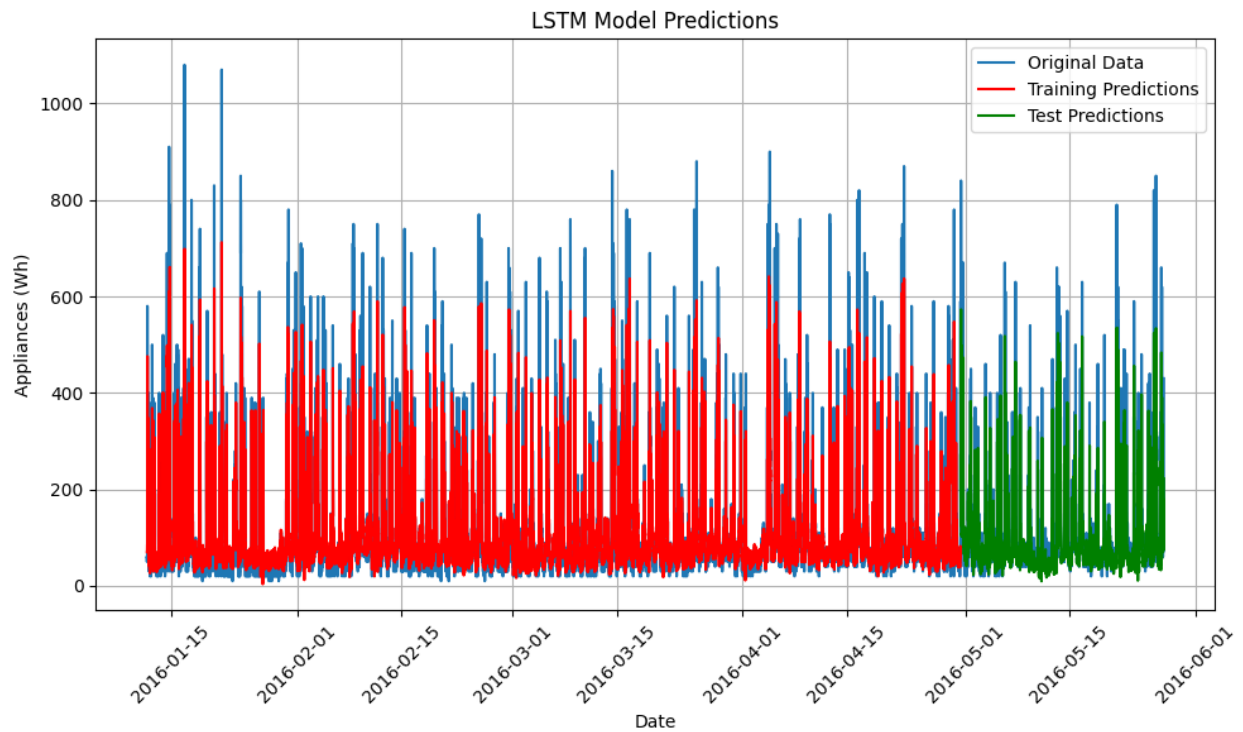
The **ReLU** activation function was used for non-linearity, and **Adam** optimizer was chosen for its efficiency in training deep networks. The model was compiled with **Mean Squared Error (MSE)** as the loss function, suitable for continuous value regression tasks like energy prediction.

Results

ARIMA Model Predictions



LSTM Model Predictions



Model Optimization

Optimization Techniques

To improve the model, time-based **cyclical features** (sine/cosine transformations of hour and weekday) and **lag features** were engineered to capture seasonal and temporal patterns. Additionally, **hyperparameter tuning** was performed using **Keras Tuner** to optimize LSTM units, dropout rate, and learning rate.

After optimization, the LSTM model achieved:

- MAE: 27.98
- RMSE: 63.28
- MAPE: 23.71%
- R² Score: 0.52

These metrics show a notable improvement in prediction accuracy and generalization compared to initial untuned models.

Challenges and Solutions

This was my first time-series analysis project. I lacked a clear understanding of the process at first but overcame this by referring to online resources and revisiting concepts from previous prediction projects. I had to work with new models like ARIMA and LSTM, which I wasn't previously familiar with. While model performance was modest, this project laid a strong foundation, and with more time and exploration, I plan to further improve accuracy in future work.

Conclusion

Model scores of ARIMA and LSTM are very close, placing them in the top tier. Following closely are Gradient Boosting and Random Forest.

Model selection and parameter tuning are complex but crucial steps in time-series forecasting. Given the 10-minute interval data over 4+ months, both **ARIMA** and **LSTM** proved to be effective. ARIMA, with basic parameter tuning, achieved reasonable performance due to its ability to leverage time-series characteristics.

While traditional machine learning models like **Random Forest** and **Gradient Boosting** can be applied, they are generally more complex to implement in time-series contexts and yielded less optimal results in this case.

References

- Time Series Forecasting:

[Time Series Forecasting with LSTM Neural Networks \(TensorFlow Tutorial\)](#)

- Feature Engineering for Time Series:

[Kaggle - Time Series Feature Engineering](#)