## Capstone Project 1: Data Visualization

**Appliances Energy Prediction –**

Business Problem Description – Dataset contains the house temperature and humidity conditions were monitored with a ZigBee wireless sensor network. As per the description on UCI website, each wireless node transmitted the temperature and humidity conditions around 3.3 min, then, the wireless data was averaged for 10 minutes periods. The energy data was logged every 10 minutes with m-bus energy meters. Combining this data with the weather data based on the date time columns

Link for the Jupyter file
- https://github.com/arijitsinha80/Springboard/blob/master/Project/Capstoneproject_Phase2.ipynb

From the Data Wrangling activity, we created the **input.csv** as the final dataset. This has 19735 observations and 30 attributes.

Divide the data in dimension wise to explore from the input dataset -

```
# Temperature sensors columns
temp_cols = ["T1", "T2", "T3", "T4", "T5", "T6", "T7", "T8", "T9"]

# Humidity sensors columns
hum_cols = ["RH_1", "RH_2", "RH_3", "RH_4", "RH_5", "RH_6", "RH_7", "RH_8"
, "RH_9"]

# Weather data columns
wth_cols = ["T_out", "Tdewpoint", "RH_out", "Press_mm_hg", "Windspeed", "V
isibility"]

# Target column
tgt = ["Appliances"]
```
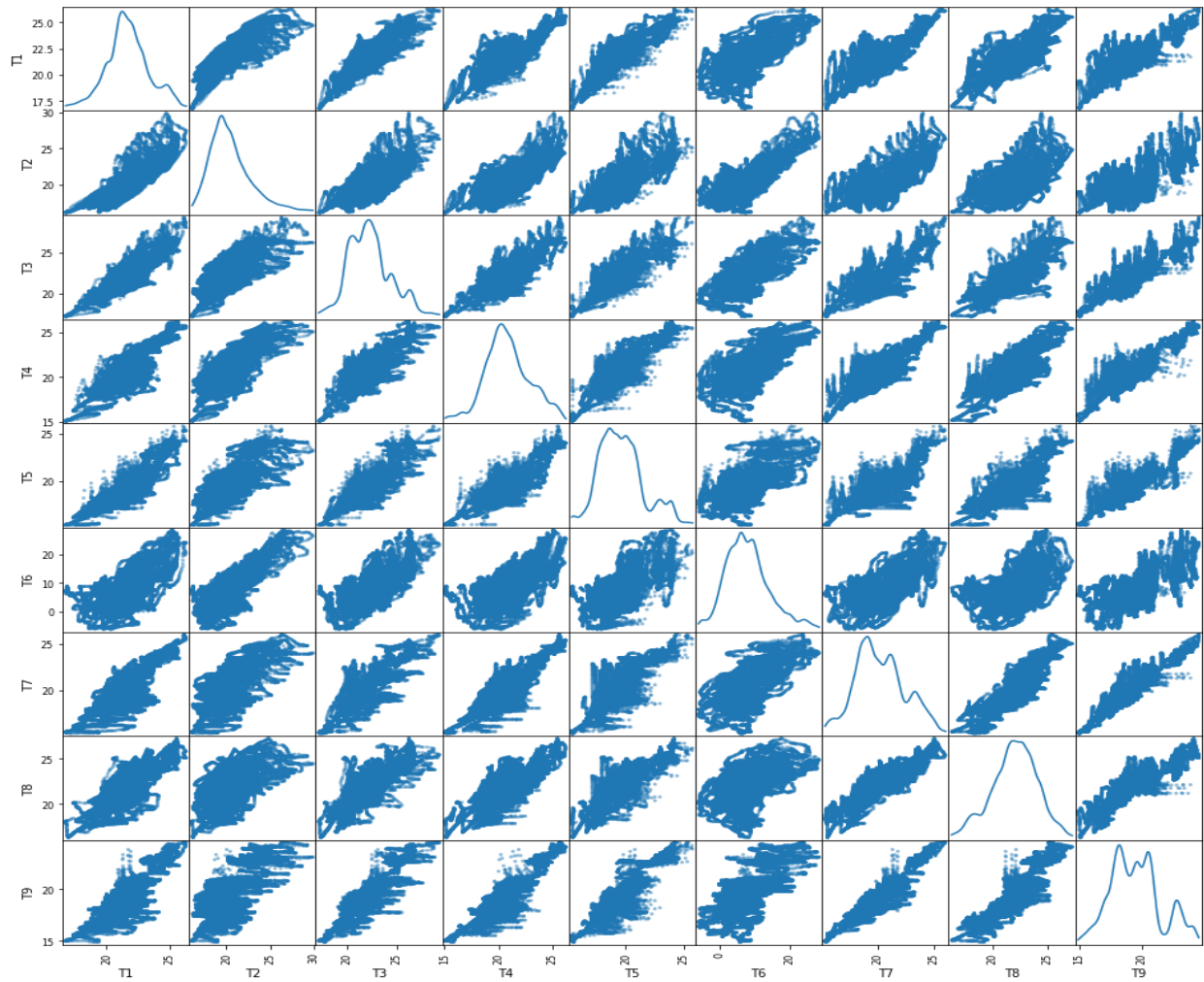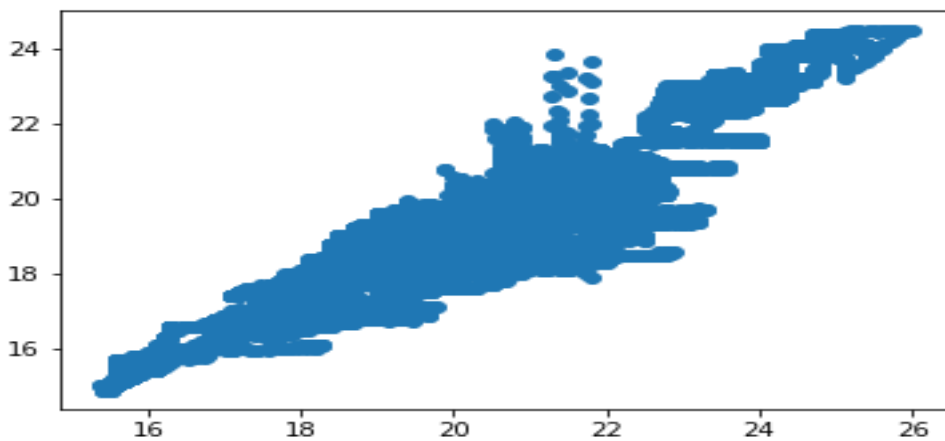
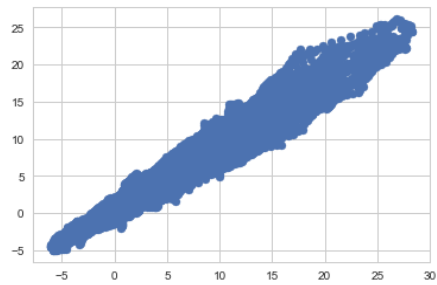From the above dimensions, we will start to explore data for each –

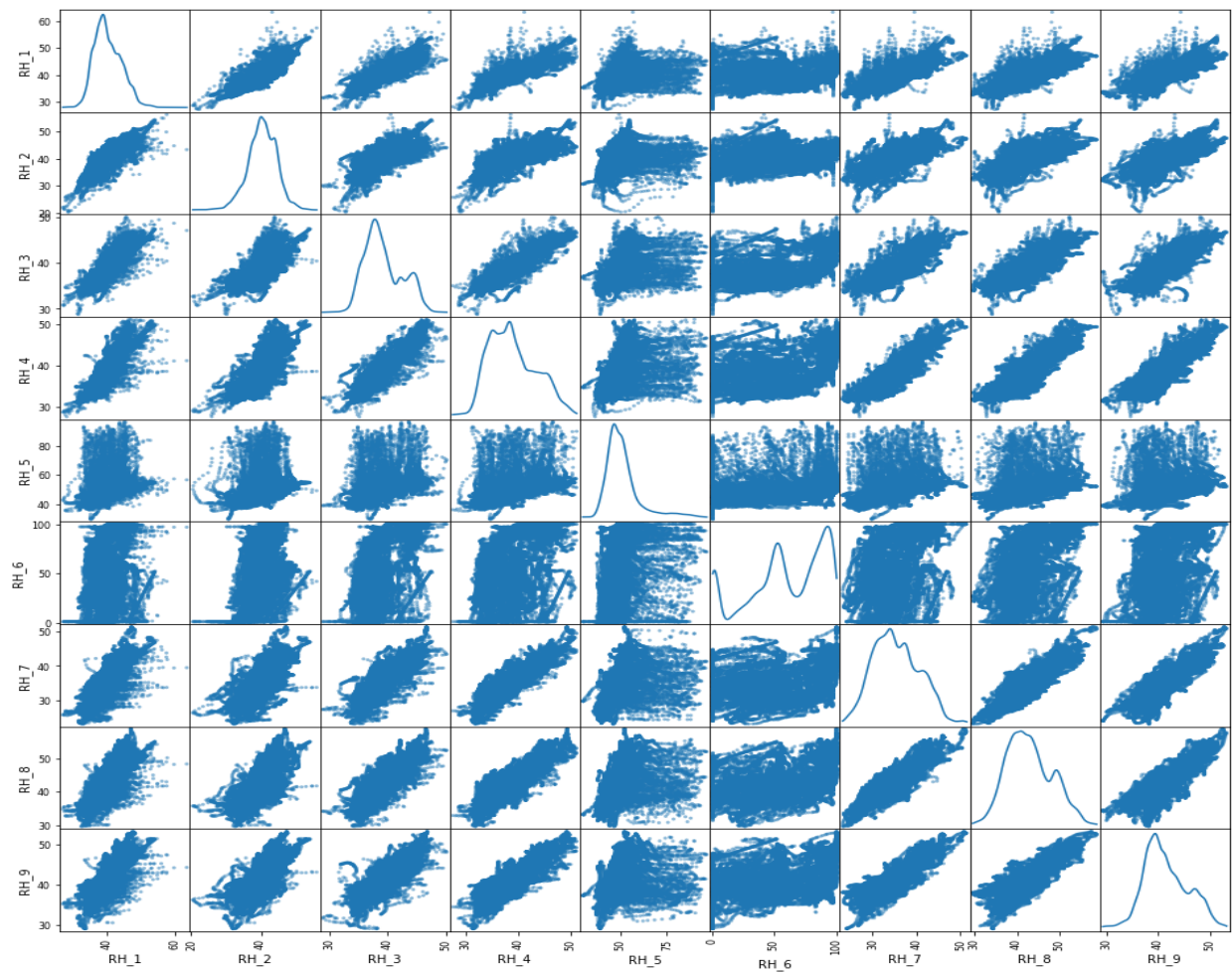Plot the scatter matrix for Temperature attributes using method **" diagonal="kde"**

- From the above figure, we can see that there is some linear relation between T7 and T9. Others are having the shape but are not exactly linear.

- There is a relation between these two attributes but also have some outliers

T6 and T_out is highly correlated, T6 is from the outside the house reading and T_out is the data collected from weather's site
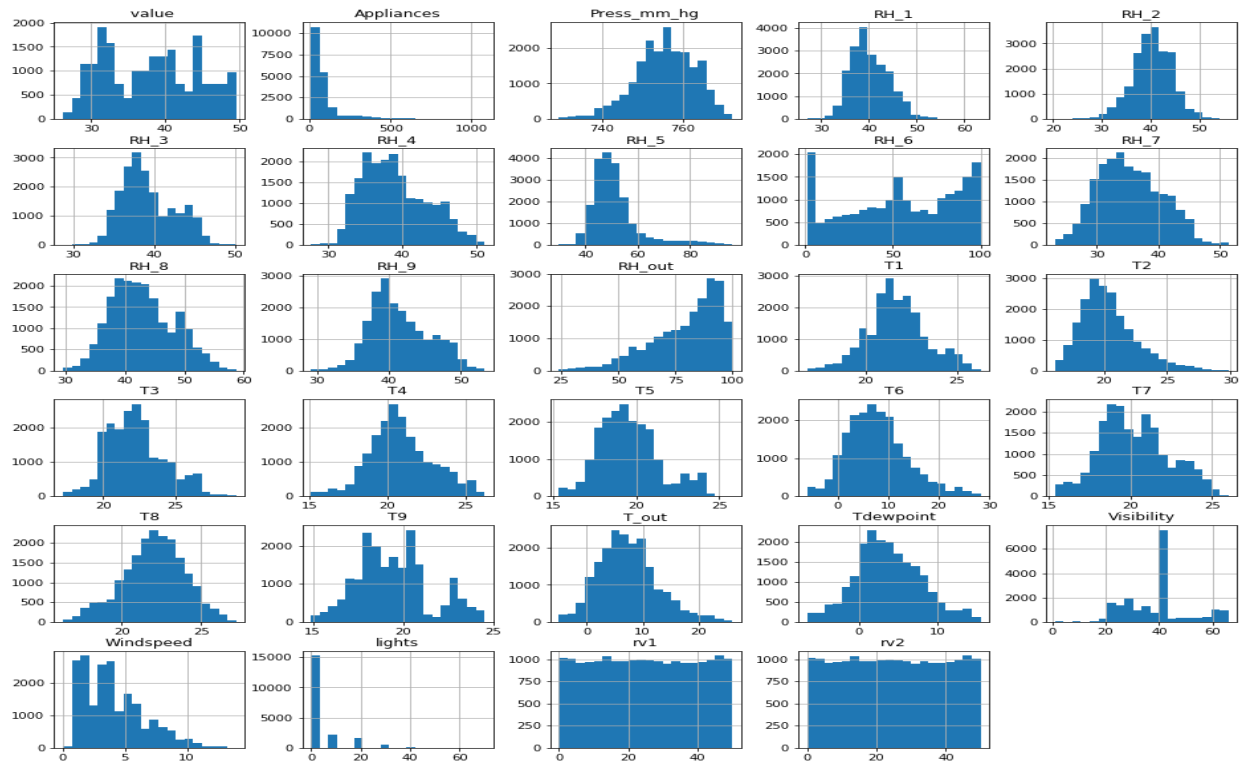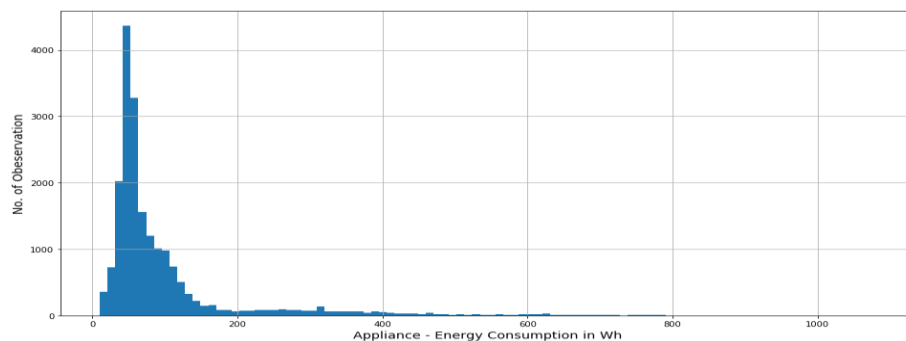


Explore the data using the Weather Dimension -



- There doesn't seems to be having any linearity between any of the attributes.
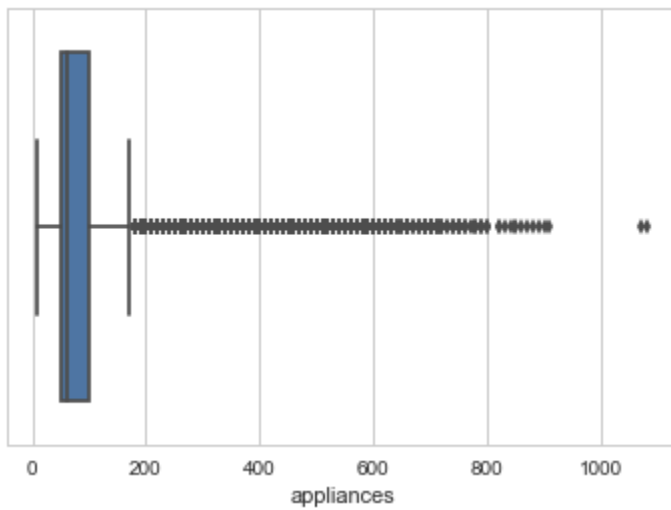
Lets explore the distribution using the histogram –



- All humidity values are almost having normal distribution except RH_6 and RH_out. In other words the reading from inside the home is having normal distribution.

- All temperature readings follow a Normal distribution except for T9.

- Visibility, Windspeed and Appliances are having skewed data.

- Rv1 and Rv2 are random variables and doesn't seems to be contributing

On the Target Attribute – Appliance, the below histograms is rightly skewed and most of the data is with 200 KWh.
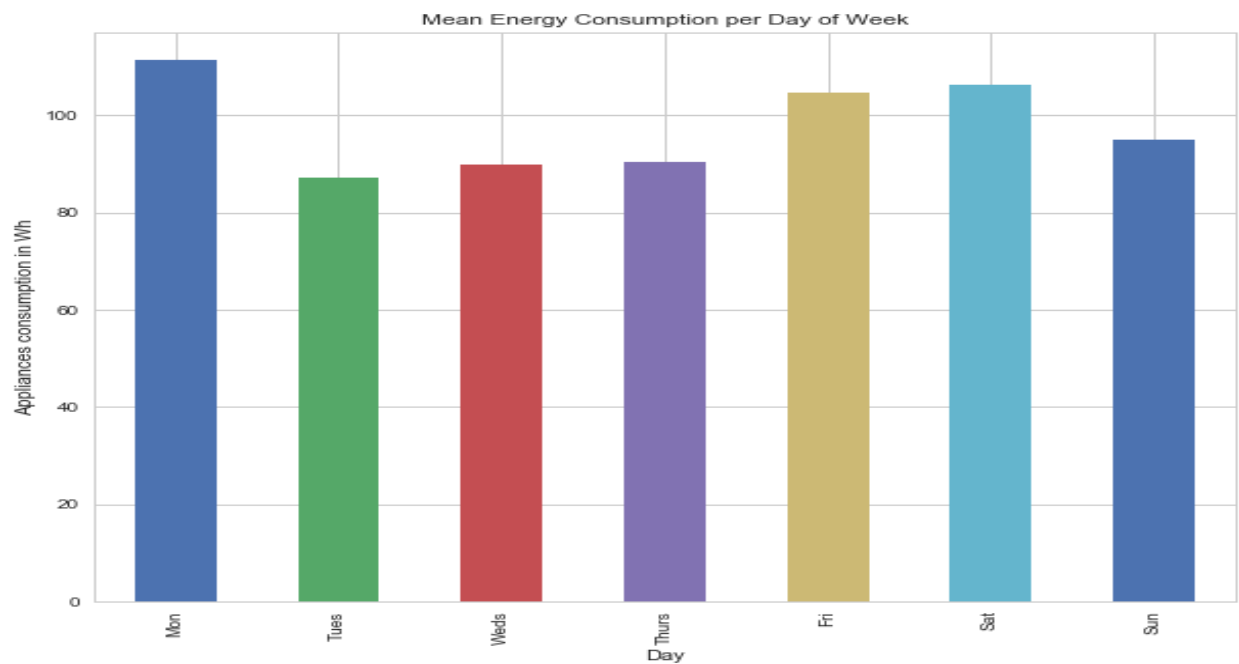
Target variable, Appliances is highly right skewed.
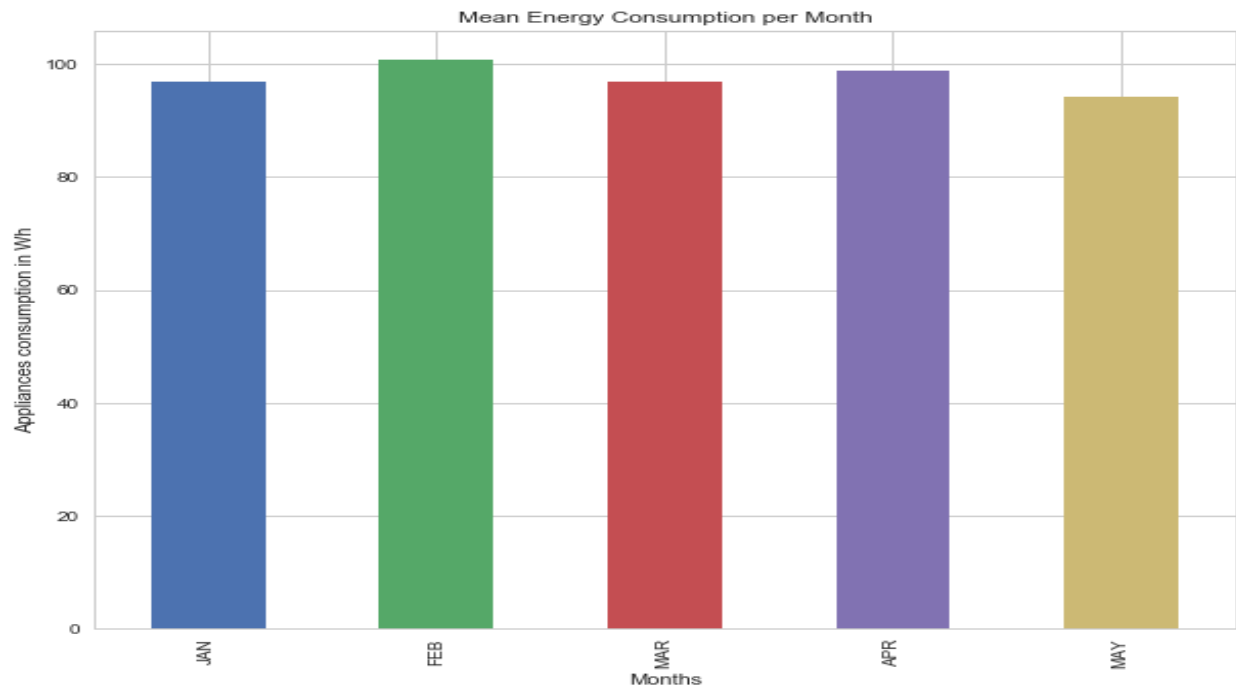Alternatively exploring using Boxplot – on Appliance Attribute



- Percentage of dataset in range of 0-200 KWh is 90.291%

Using the date attributes, created new columns for Month and Weeks using the **datetimestamp** method
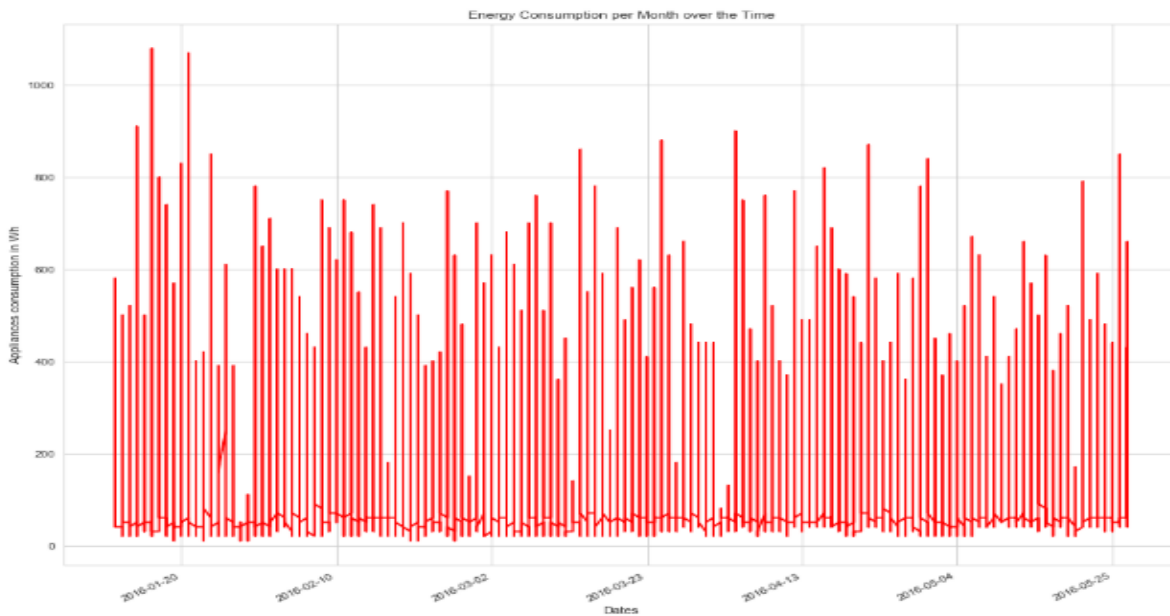
With taking the average on week – Monday the usage has been higher, followed by Saturday and Friday.
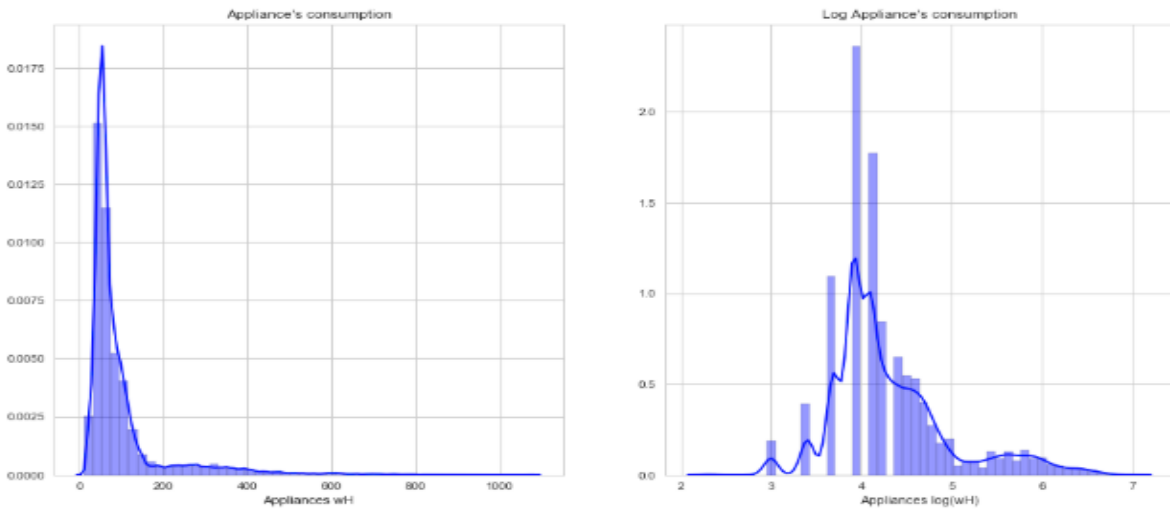
For Monthly Average – On a average, February and April the consumption has been more than other months.



For Day wise consumption – plotting this date wise, energy consumption, In January month there were 2 days when the consumption was more than 1000 KWh.
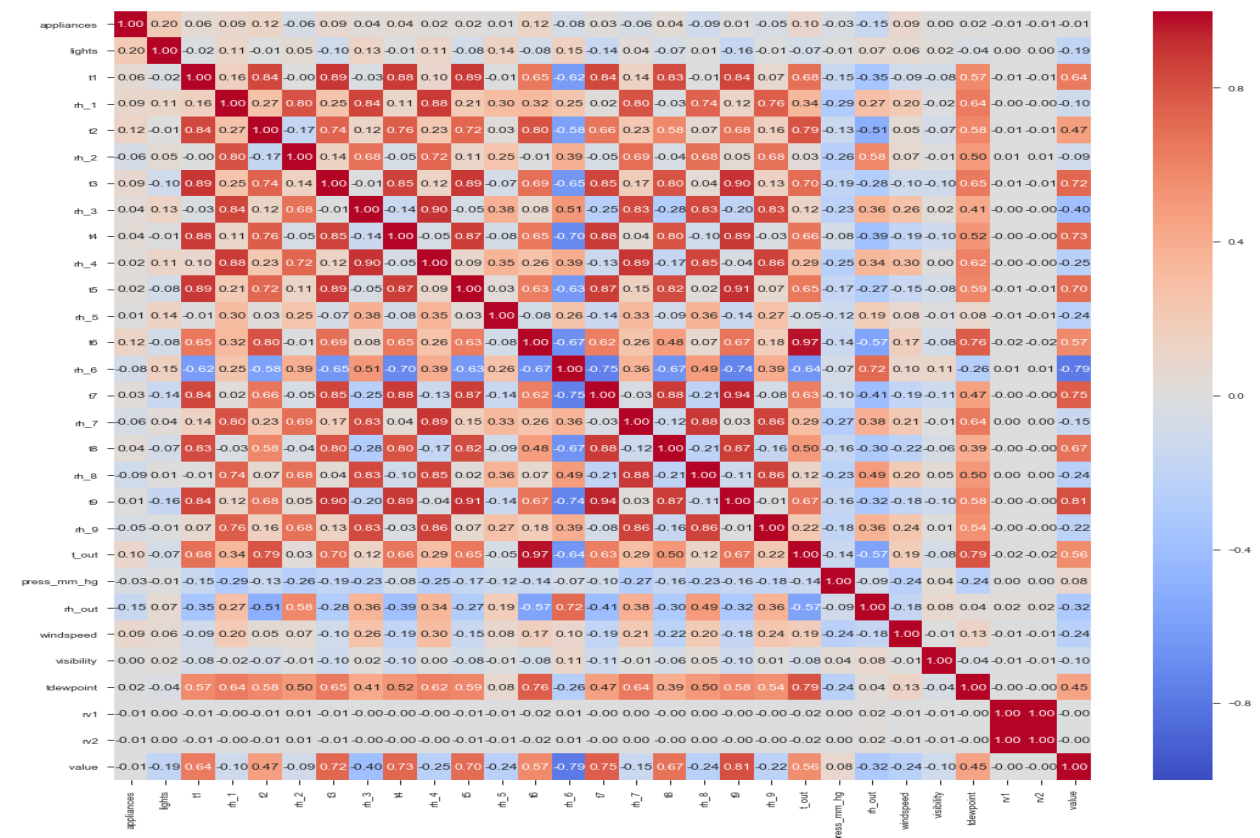
Since the Appliance data captured were rightly skewed, converting the column to Log values to see if it has the normal distribution.
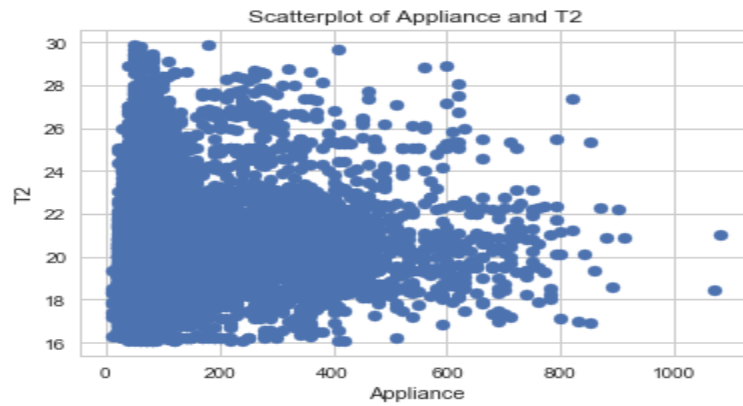


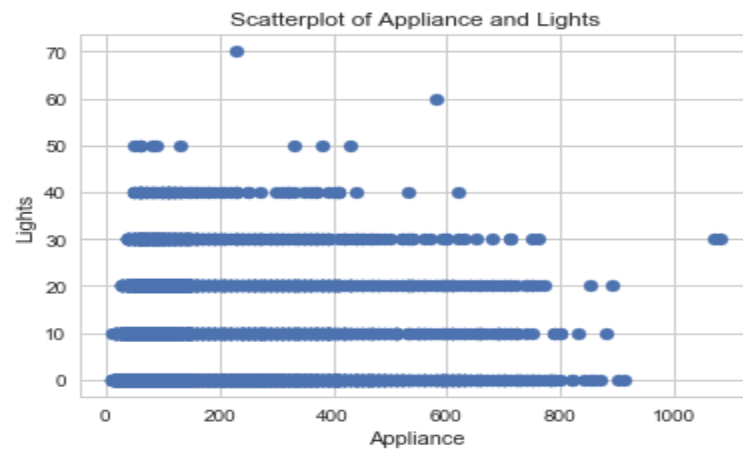Let's explore the Correlation plot –
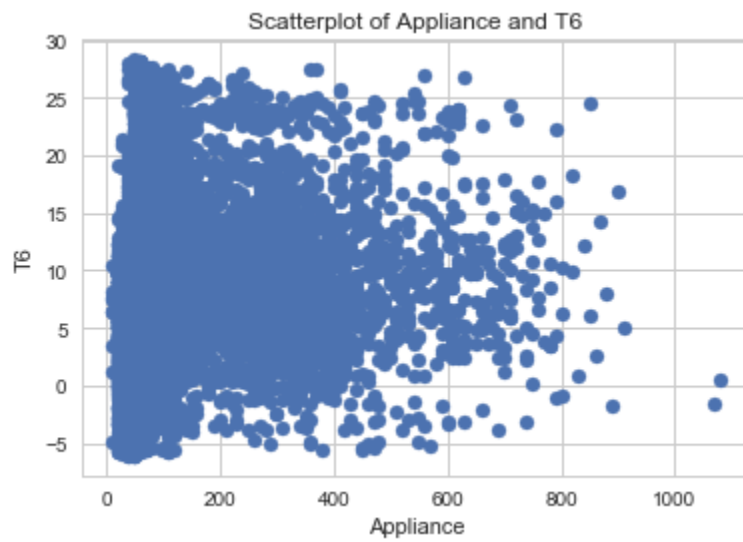
With Appliance attribute –

Scatterplot between appliances and t2


Scatterplot of Appliance and T2

Scatter plot between Appliance and Lights


Scatterplot of Appliance and Lights

Scatterplot between Appliance and T6
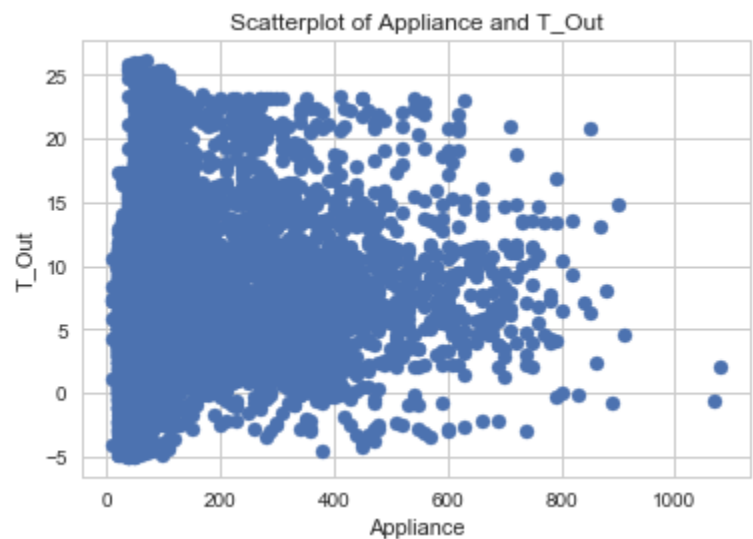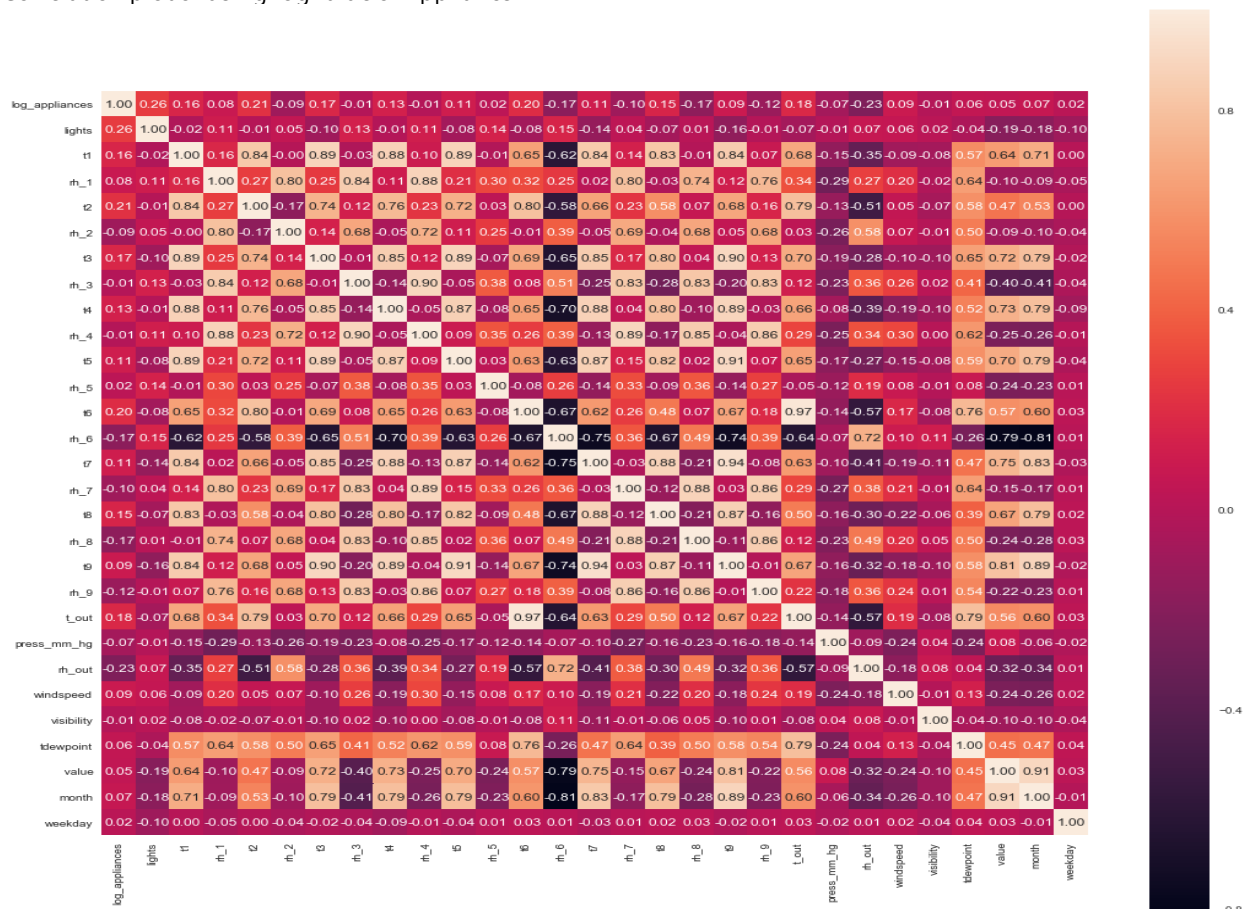

Scatterplot of Appliance and T6

Scatterplot between Appliance and T_out
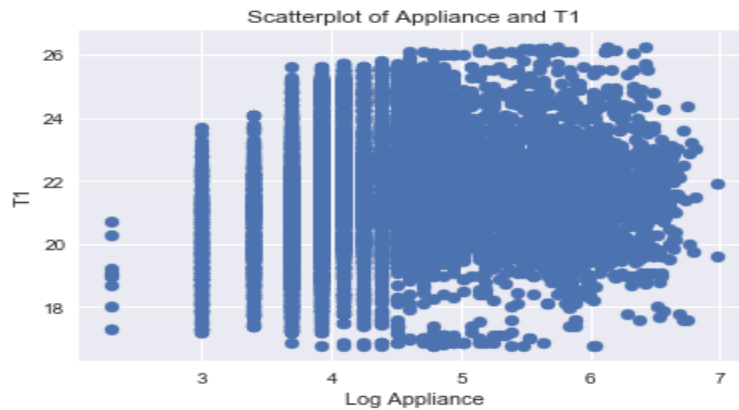


Scatterplot of Appliance and T_Out
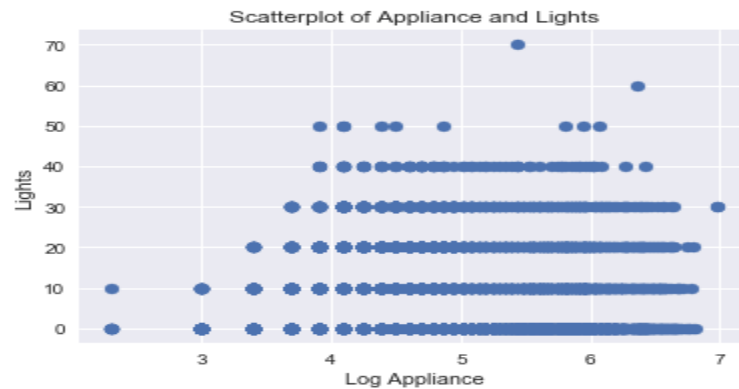
Correlation plot of using Log value of Appliance -

With Log Appliance
- The most correlated features with energy consumption(log_appliances) are: lights=0.26, t6=0.20, t2=0.22, t3 = 0.17,t_out = 0.18, rh_out = -0.23, rh_8 = -0.17, rh_6 = -0.17, windspeed = 0.09.

- In a linear regression problem only linear independent variables can be be used as features to explain energy consumption otherwise we will have multicollinearity issues.

Scatter plot of log_appliance and t1


Scatterplot of Appliance and T1

Scatterplot between log_appliance and lights


Scatterplot of Appliance and Lights

Scatter plot between log appliance and t_out


Scatterplot of Appliance and T_out