

## Capstone Project 1: Stats

### 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\\_stats.ipynb](https://github.com/arijitsinha80/Springboard/blob/master/Project/Capstoneproject_Phase2_stats.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"]

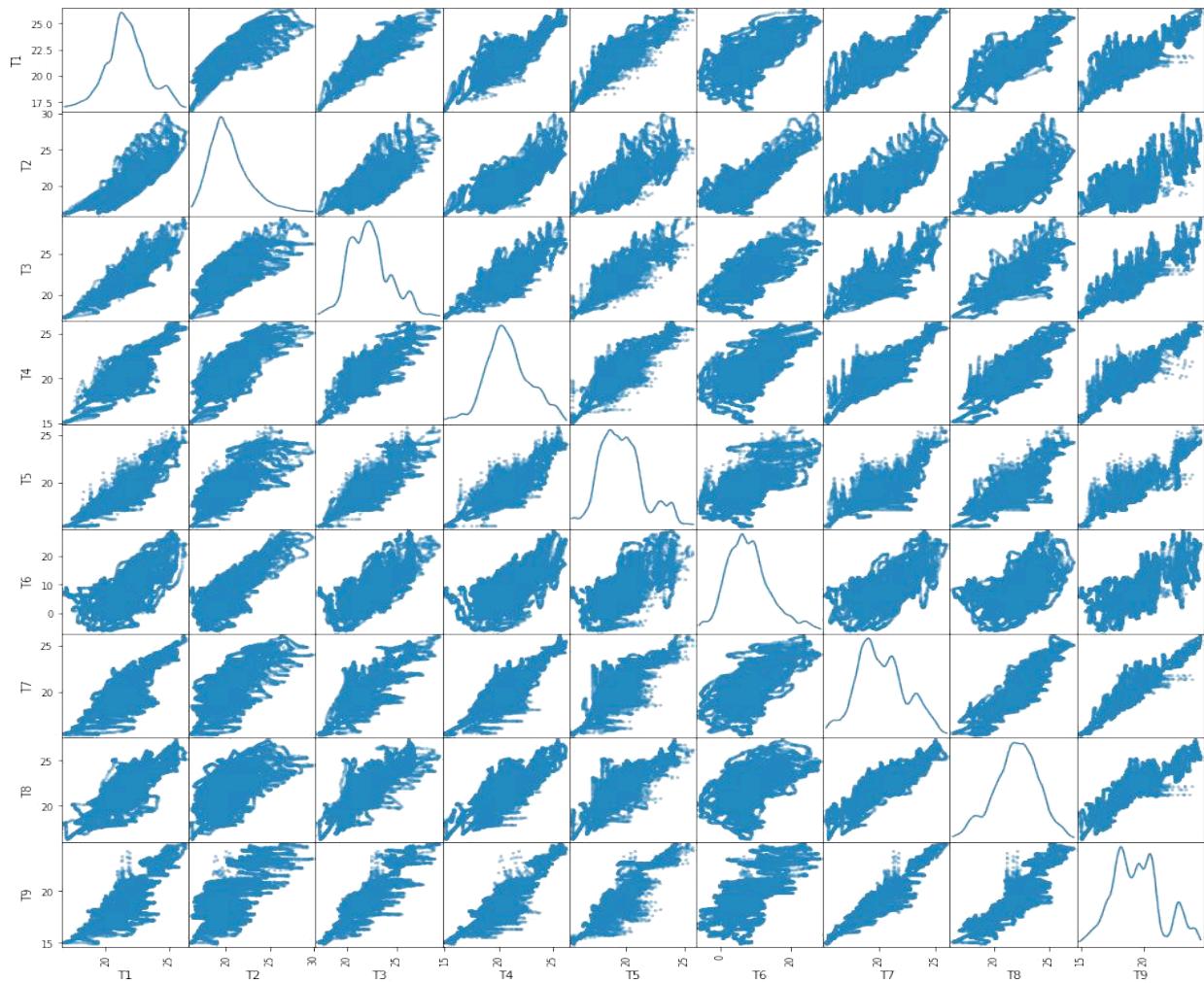
# Humidity sensors columns
hum_cols = ["RH_1", "RH_2", "RH_3", "RH_4", "RH_5", "RH_6", "RH_7", "RH_8",
, "RH_9"]

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

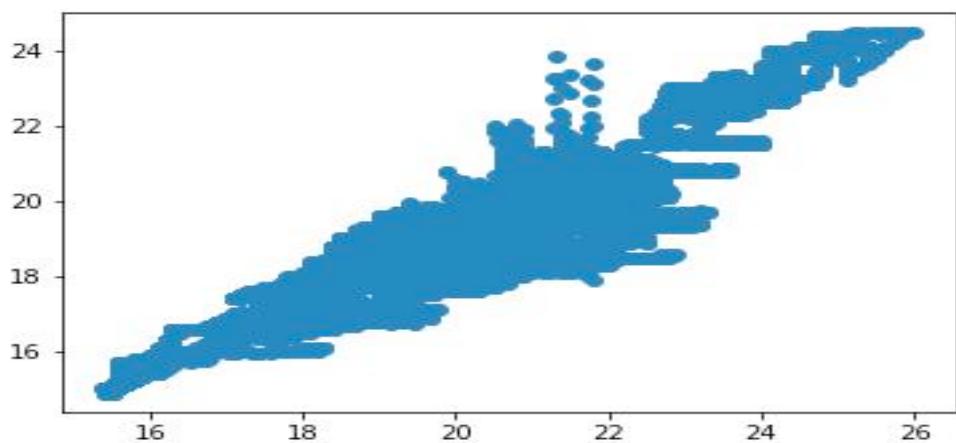
# Target column
tgt = ["Appliances"]
```

From the above dimensions, we will start to explore data for each – Plot the scatter matrix for Temperature attributes using method “ **diagonal="kde"** ”

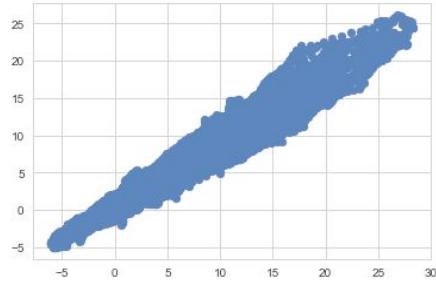
Plot the scatter matrix for Temperature attributes using method “ **diagonal="kde"** ”



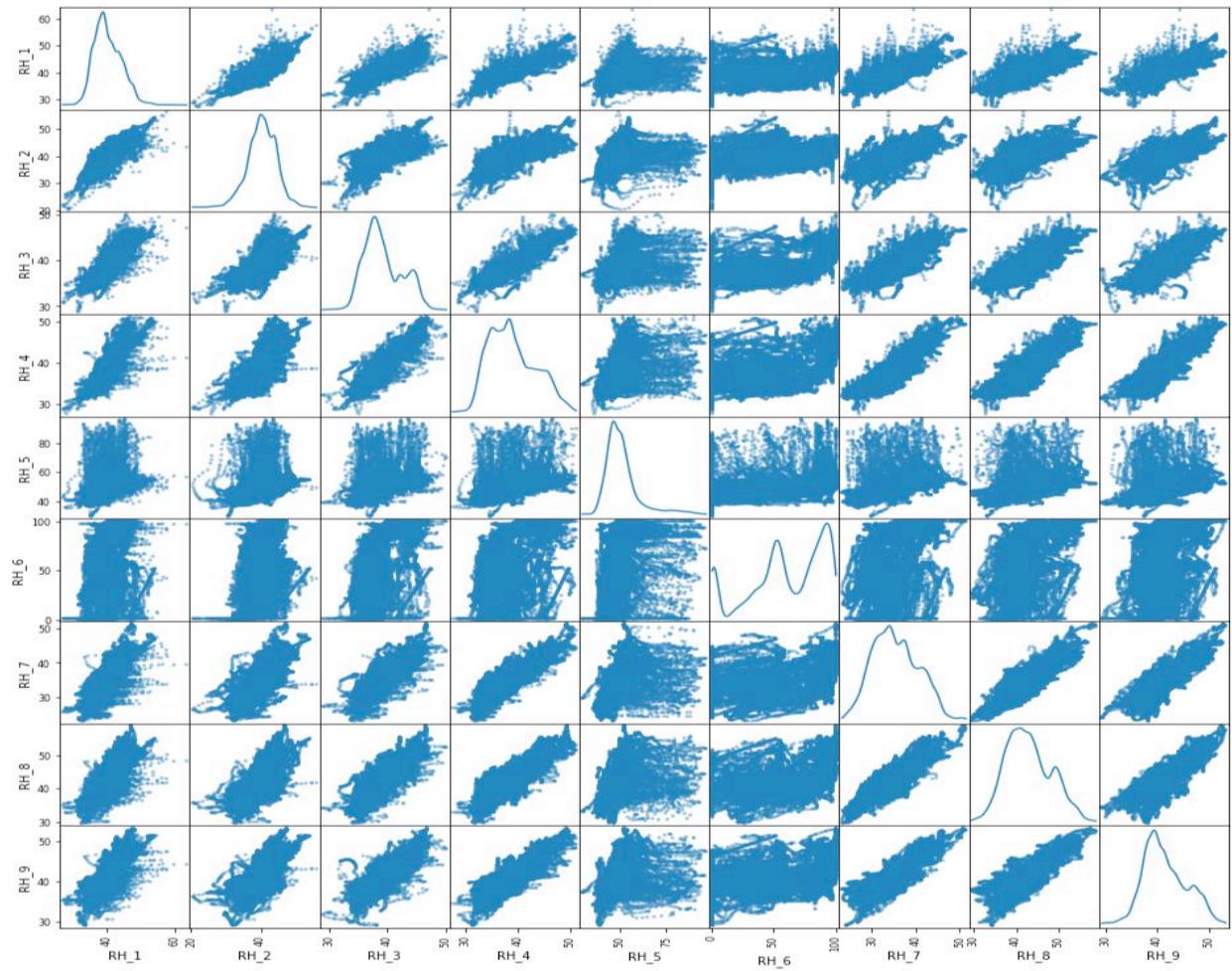
- From the above figure, we can see that there is some linear relation between T7 and T9. Others are having the shape but are not exactly linear.
- There is a relation between these two attributes but also have some outliers



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

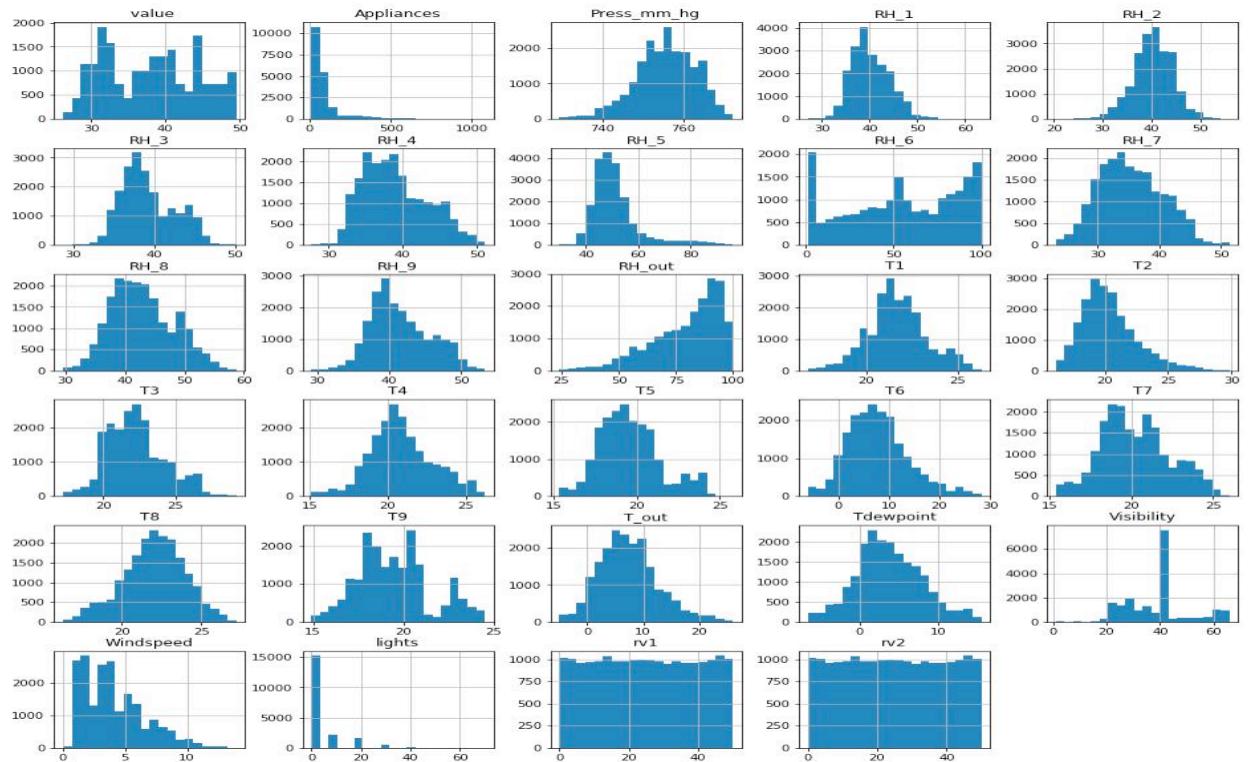


Explore the data using the Weather Dimension -



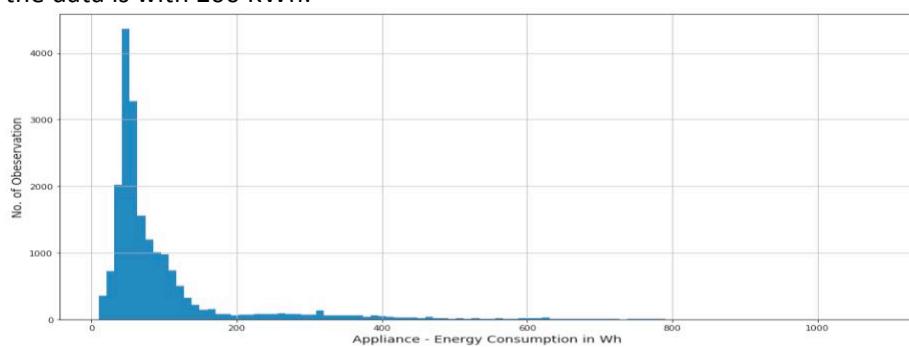
- There doesn't seem to be having any linearity between any of the attributes.

Lets explore the distribution using the histogram –

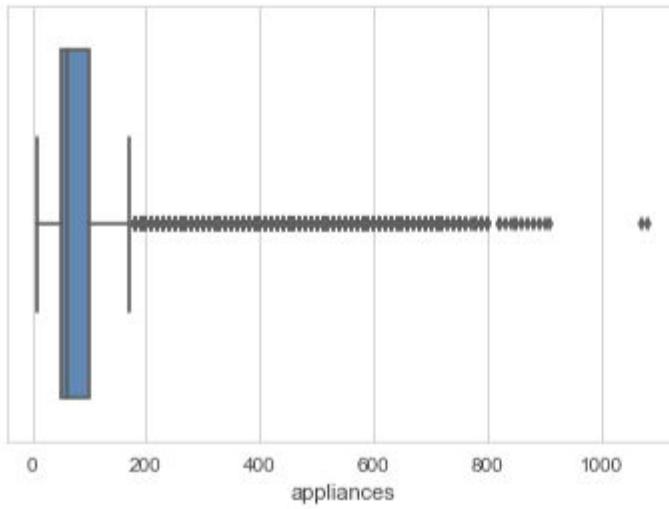


- All humidity values are almost having normal distribution except RH\_6 and RH\_out. In other words the reading from inside the home is having normal distribution.
- All temperature readings follow a Normal distribution except for T9.
- Visibility, Windspeed and Appliances are having skewed data.
- 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.



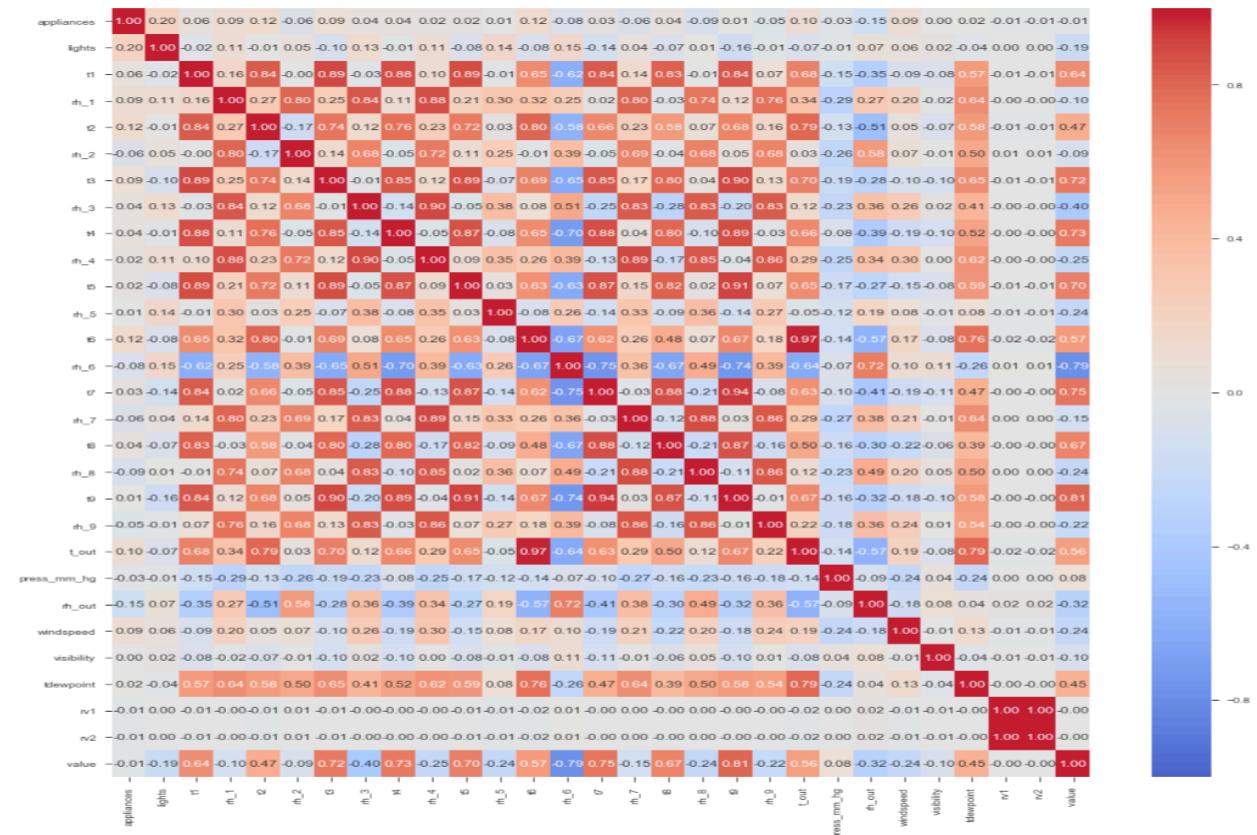
Target variable, Appliances is highly right skewed.  
Alternatively exploring using Boxplot – on Appliance Attribute



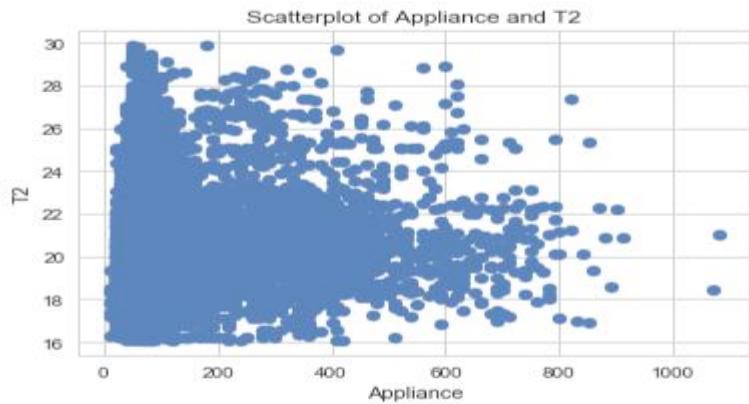
- Percentage of dataset in range of 0-200 KWh is 90.291%

Let's explore the Correlation plot –

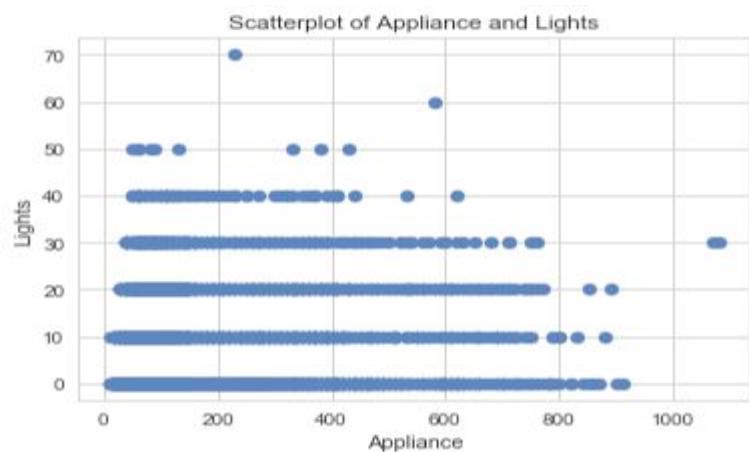
With Appliance attribute –



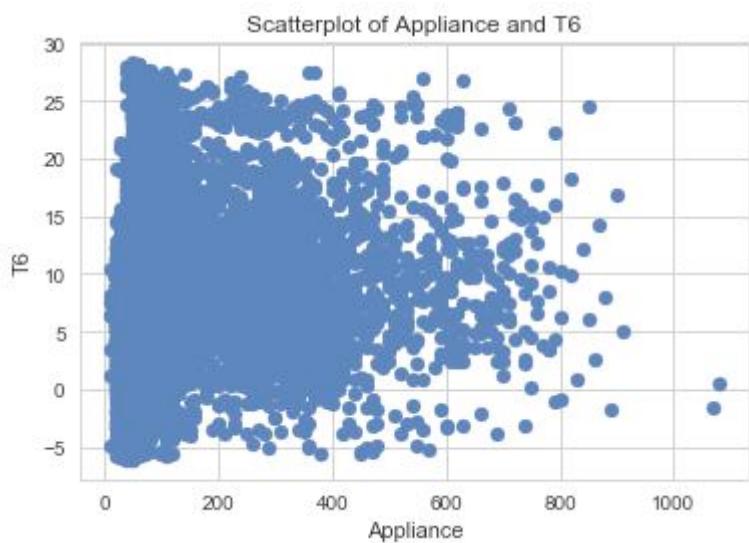
Scatterplot between appliances and t2



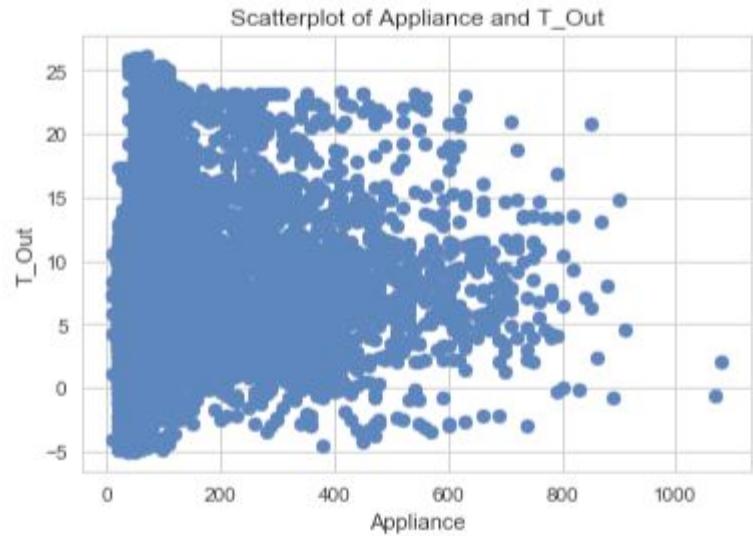
Scatter plot between Appliance and Lights



Scatterplot between Appliance and T6

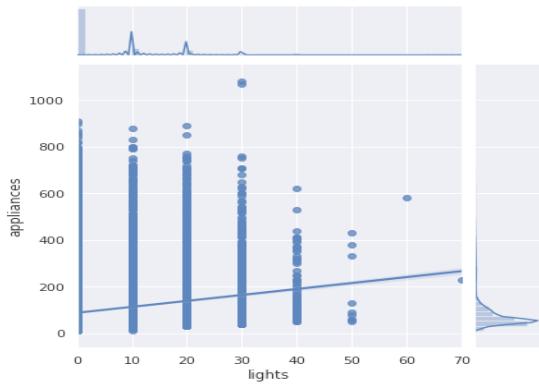


### Scatterplot between Appliance and T\_out

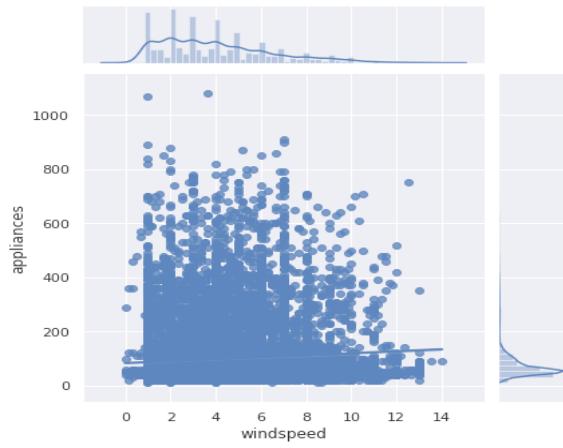


**Variables that are particularly significant in terms of predicting Appliance Energy Consumption based on the correlation matrix –**

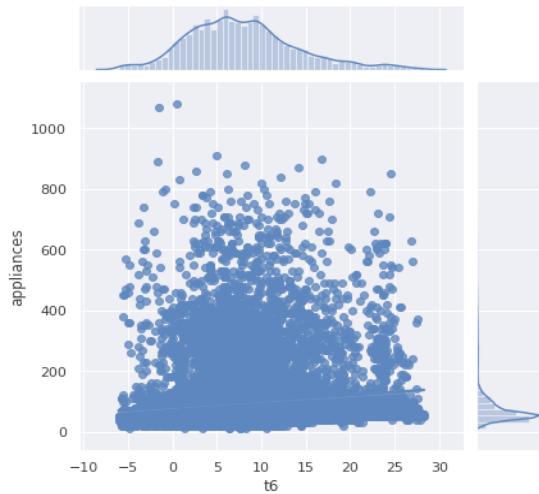
- Between Appliance and Lights



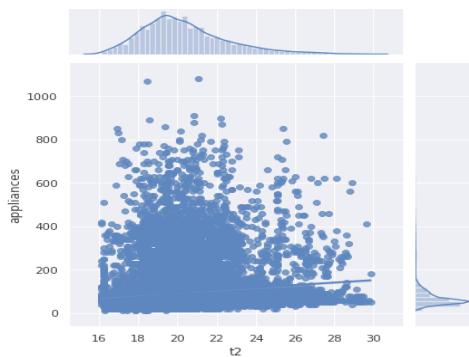
- Between Appliance and Windspeed



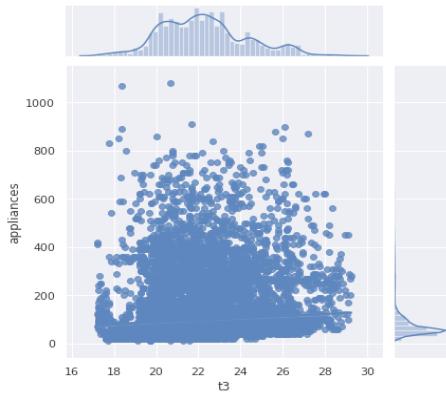
- Between Appliance and T6



- Between Appliance and T2



- Between Appliance and T3



### Calculate the Correlation between Temperature features -

- Correlation between T9 and T1 pearson 0.84 0.00 None
- Correlation between T9 and T2 pearson 0.68 0.00 None
- Correlation between T9 and T3 pearson 0.90 0.00 None
- Correlation between T9 and T4 pearson 0.89 0.00 None
- Correlation between T9 and T5 pearson 0.91 0.00 None
- Correlation between T9 and T6 pearson 0.67 0.00 None
- Correlation between T9 and T7 pearson 0.94 0.00 None
- Correlation between T9 and T8 pearson 0.87 0.00 None

### Check, if the Temperature, Humidity and Weather features influences Appliance – Using OLS Summary

1. Coefficient table (middle table). We can interpret the t3 coefficient (4.3471) by first noticing that the p-value (under  $P>|t|$ ) is so small, basically zero. This means that the t3 is a statistically significant predictor of appliance energy consumption.

The regression coefficient for t3 of 4.3471 means that on average, each additional t3 temperature is associated with an increase the appliance energy consumption

The confidence interval gives us a range of plausible values for this average change, about (3.637 and 5.058)

$R^2$  is only 0.007, hence t3 doesn't contribute much on the variance. F-Statistic The F-Statistic is 143.8 and the probability for this statistic is 5.09e-33, which is close to 0. We can safely reject the null hypothesis, indicating that at least one  $\beta$  coefficient is nonzero.

```

OLS Regression Results
=====
Dep. Variable: appliances R-squared: 0.007
Model: OLS Adj. R-squared: 0.007
Method: Least Squares F-statistic: 143.8
Date: Mon, 18 May 2020 Prob (F-statistic): 5.09e-33
Time: 21:12:10 Log-Likelihood: -1.1931e+05
No. Observations: 19735 AIC: 2.386e+05
Df Residuals: 19733 BIC: 2.386e+05
Df Model: 1
Covariance Type: nonrobust

coef std err t P>|t| [0.025 0.975]
Intercept 0.8955 8.105 0.110 0.912 -14.990 16.781
t3 4.3471 0.362 11.992 0.000 3.637 5.058

Omnibus: 14099.091 Durbin-Watson: 0.498
Prob(Omnibus): 0.000 Jarque-Bera (JB): 196052.616
Skew: 3.410 Prob(JB): 0.00
Kurtosis: 16.854 Cond. No. 250.

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

2. Coefficient table (middle table). We can interpret the t3+t6 coefficient (0.4119, 1.8871) by first noticing that the p-value (under  $P>|t|$ ) is so small, basically zero. This means that the t6 is a statistical significant predictor of appliance energy consumption.

The regression coefficient for t6 of 1.8871, means that on average, each additional t6 temperature is associated with an increase the appliance energy consumption

The confidence interval gives us a range of plausible values for this average change, about (1.566 and 2.208)

$R^2$  is only 0.014, hence t3 and t6 doesn't contribute much on the variance. F-Statistic The F-Statistic is 138.8 and the probability for this statistic is 1.39e-60, which is close to 0. We can safely reject the null hypothesis, indicating that at least one  $\beta$  coefficient is nonzero.

```

OLS Regression Results
=====
Dep. Variable: appliances R-squared: 0.014
Model: OLS Adj. R-squared: 0.014
Method: Least Squares F-statistic: 138.8
Date: Mon, 18 May 2020 Prob (F-statistic): 1.39e-60
Time: 21:12:21 Log-Likelihood: -1.1924e+05
No. Observations: 19735 AIC: 2.385e+05
Df Residuals: 19732 BIC: 2.385e+05
Df Model: 2
Covariance Type: nonrobust

coef std err t P>|t| [0.025 0.975]
Intercept 73.5949 10.249 7.181 0.000 53.506 93.684
t3 0.4119 0.497 0.828 0.407 -0.563 1.386
t6 1.8871 0.164 11.525 0.000 1.566 2.208

Omnibus: 14117.484 Durbin-Watson: 0.500
Prob(Omnibus): 0.000 Jarque-Bera (JB): 197909.447
Skew: 3.412 Prob(JB): 0.00
Kurtosis: 16.932 Cond. No. 339.

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

3. Coefficient table (middle table). We can interpret the t3+t6+rh\_out coefficient (1.8057, 0.3079, -0.9076) by first noticing that the p-value (under  $P>|t|$ ) is so small, basically zero. This means that the t6 is a statistical significant predictor of appliance energy consumption.

The regression coefficient for rh\_out of -0.9076, means that on average, each additional t6 temperature is associated with an decrease the appliance energy consumption

The confidence interval gives us a range of plausible values for this average change, about (-1.025 and -0.790)

$R^2$  is only 0.025 better than previous, hence t3, t6 and rh\_out doesn't contribute much on the variance. F-Statistic The F-Statistic is 170.3 and the probability for this statistic is 4.96e-109, which is close to 0. We can safely reject the null hypothesis, indicating that at least one  $\beta$  coefficient is nonzero.

```
OLS Regression Results
=====
Dep. Variable: appliances R-squared: 0.025
Model: OLS Adj. R-squared: 0.025
Method: Least Squares F-statistic: 170.3
Date: Mon, 18 May 2020 Prob (F-statistic): 4.96e-109
Time: 21:12:55 Log-Likelihood: -1.1913e+05
No. Observations: 19735 AIC: 2.383e+05
Df Residuals: 19731 BIC: 2.383e+05
Df Model: 3
Covariance Type: nonrobust
=====
      coef  std err      t      P>|t|      [0.025      0.975]
Intercept  127.4360   10.790    11.810     0.000    106.286    148.586
t3          1.8057    0.503     3.592     0.000     0.820     2.791
t6          0.3079    0.193     1.593     0.111    -0.071     0.687
rh_out     -0.9076    0.060    -15.173     0.000    -1.025    -0.790
=====
Omnibus: 14135.525 Durbin-Watson: 0.507
Prob(Omnibus): 0.000 Jarque-Bera (JB): 199836.330
Skew: 3.415 Prob(JB): 0.00
Kurtosis: 17.014 Cond. No. 1.26e+03
=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.26e+03. This might indicate that there are
strong multicollinearity or other numerical problems.
```

4. Coefficient table (middle table). We can interpret the  $t_1+t_2+t_3+t_4+t_5+t_6+t_7+t_8+rh_{-1}+rh_{-2}+windspeed$ , coefficient (9.0446, -25.6614, 17.7293, -1.4768, -7.3830, -7.5650, 1.0356, -6.2685, 9.4475, 20.0347, -20.3286, 1.6784) by first noticing that the p-value (under  $P>|t|$ ) is so small, basically zero. This means that the  $t_6$  is a statistically significant predictor of appliance energy consumption.

The confidence interval of t3 gives us a range of plausible values for this average change, about (15.814 and 19.644)

$R^2$  is only 0.098 better than previous, F-Statistic The F-Statistic is 194.5 and the probability for this statistic is 0. We can safely reject the null hypothesis, indicating that at least one  $\beta$  coefficient is nonzero.

```

OLS Regression Results
-----
Dep. Variable: appliances R-squared: 0.098
Model: OLS Adj. R-squared: 0.097
Method: Least Squares F-statistic: 194.5
Date: Mon, 18 May 2020 Prob (F-statistic): 0.00
Time: 21:35:34 Log-Likelihood: -1.1836e+05
No. Observations: 19735 AIC: 2.367e+05
Df Residuals: 19723 BIC: 2.368e+05
Df Model: 11
Df Robust: 11
Covariance Type: nonrobust

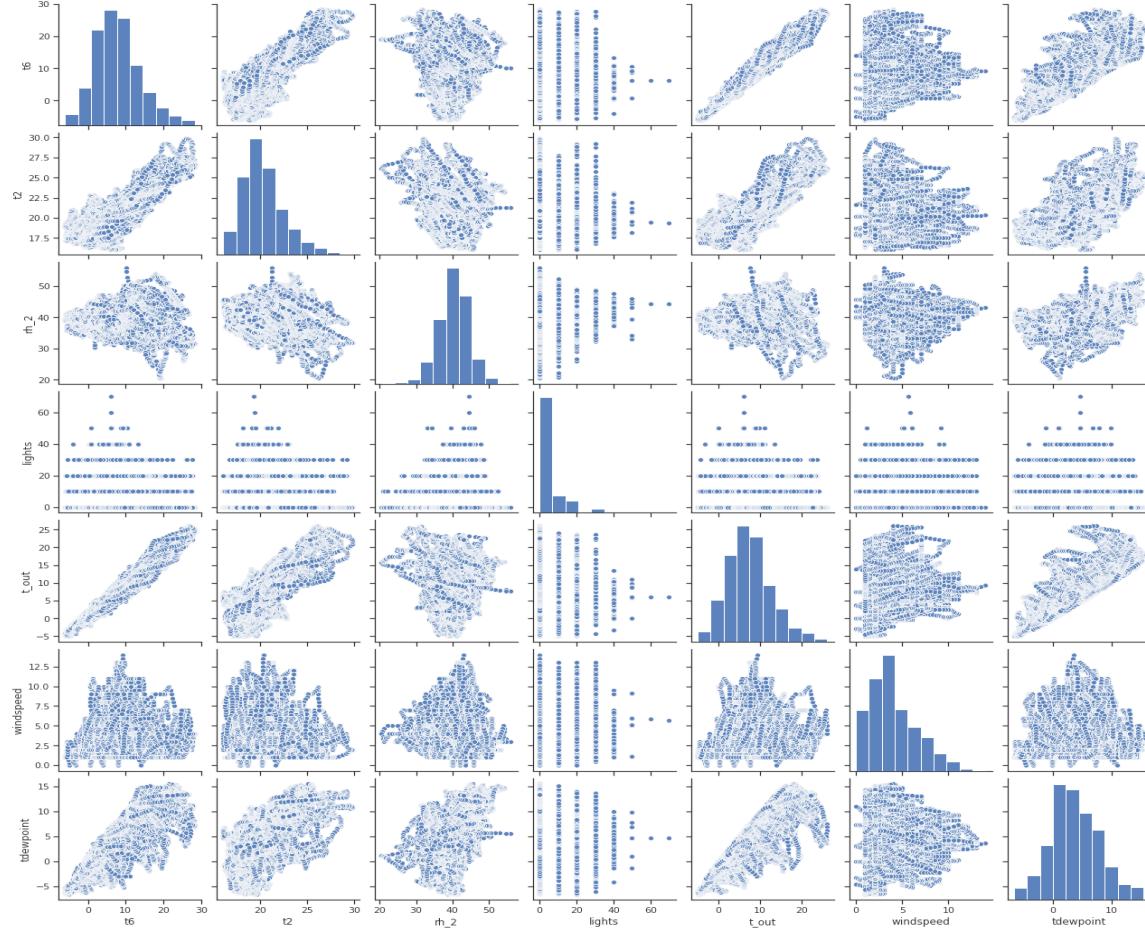
coef std err t P>|t| [0.025 0.975]
-----
Intercept 123.2455 15.477 7.963 0.000 92.909 153.582
t1 9.0446 1.764 5.127 0.000 5.587 12.502
t2 -25.6614 1.450 -17.697 0.000 -28.504 -22.819
t3 17.7293 0.977 18.148 0.000 15.814 19.644
t4 -1.4768 0.907 -1.628 0.103 -3.255 0.301
t5 -7.3830 1.065 -6.930 0.000 -9.471 -5.295
t6 1.0356 0.227 4.552 0.000 0.590 1.482
t7 -6.2685 0.971 -6.457 0.000 -8.171 -4.366
t8 9.4475 0.881 10.730 0.000 7.722 11.173
rh_1 20.0347 0.630 31.792 0.000 18.799 21.270
rh_2 -20.3286 0.630 -32.261 0.000 -21.564 -19.094
windspeed 1.6784 0.320 5.239 0.000 1.050 2.306

Omnibus: 13836.970 Durbin-Watson: 0.578
Prob(Omnibus): 0.000 Jarque-Bera (JB): 196316.384
Skew: 3.306 Prob(JB): 0.00
Kurtosis: 16.965 Cond. No. 1.80e+03

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.8e+03. This might indicate that there are
strong multicollinearity or other numerical problems.

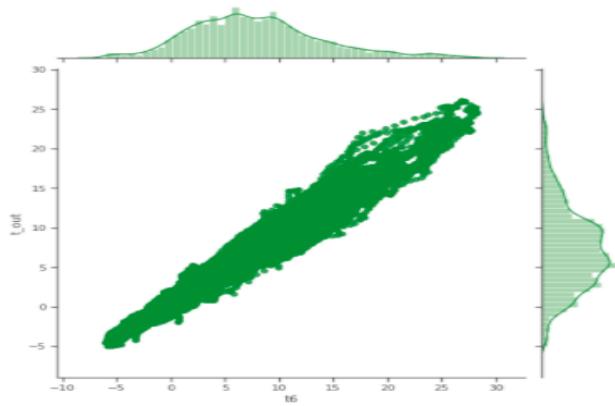
```

### Pairplot for 't6','t2','rh\_2','lights','t\_out','windspeed','tdewpoint' features for their distribution –



### **Is there a significant difference between T6 and T\_out and impact my future Prediction Models –**

With Description and plotting the **jointplot** of the two features -



**Run a Two-sided T-test with the following hypotheses:**

Null hypothesis:  $t_6 = t_{out}$

Alternate hypothesis:  $t_6 \neq t_{out}$

Upon Conducting the T-Test – received the below values - Ttest\_indResult (statistic=8.675177895656354, pvalue=4.283728402821399e-18)

Result - Given the high p-value: 4.2, hence will not reject the null hypothesis that feature t6 and t\_out almost same and redundant.

### **Major Inference –**

- Temperature feature from T1-T9 and T\_out have positive correlation with the target Appliances. For the indoor temperatures, the correlations are high as expected. Four columns have a high degree of correlation with T9 & T3,T5,T7,T8 also T6 & T\_Out has high correlation (both temperatures from outside) . Hence we can remove the T9 and T\_out from the model in next section.
- Weather attributes - Visibility, Tdewpoint, Press\_mm\_hg have low correlation values
- Humidity - There are no significantly high correlation cases for humidity sensors.
- Random variables have no role to play; hence we will remove these features from the model in next section.