

# **Capstone Project**

## **Appliances Energy Prediction**

by

**Pankaj Beldar**  
**Almabetter Trainee, Bangalore**

# Problem Statement

- Data-driven prediction of energy use of appliances .The data set is at 10 min for about 4.5 months. The house temperature and humidity conditions were monitored with a ZigBee wireless sensor network. 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. Weather from the nearest airport weather station (Chievres Airport, Belgium) was downloaded from a public data set from Reliable Prognosis (rp5.ru) and merged together with the experimental data sets using the date and time column. Two random variables have been included in the data set for testing the regression models. and to filter out non-predictive attributes (parameters).

# Introduction

- The understanding of the appliances energy use in buildings has been the subject of numerous research studies , since appliances represent a significant portion (between 20 and 30% of the electrical energy demand .
- Regression models for energy use can help to understand the relationships between different variables and to quantify their impact. Thus, prediction models of electrical energy consumption in buildings can be useful for a number of applications: to determine adequate sizing of photovoltaic and energy storage to diminish power flow into the grid , to detect abnormal energy use patterns , to be part of an energy management system for load control , to model predictive control applications where the loads are needed , for demand side management (DSM) and demand side response (DSR) and as an input for building performance simulation analysis

# Data Attributes

- **date:** time year-month-day  
hour:minute:second
- **Appliances:** energy use in Wh (Dependent variable)
- **lights:** energy use of light fixtures in the house in Wh (Drop this column)
- **T1,** Temperature in kitchen area, in Celsius
- **RH1:** Humidity in kitchen area, in %
- **T2,** Temperature: in living room area, in Celsius
- **RH2:**Humidity in living room area, in %
- **T3:** Temperature in laundry room area
- **RH3:** Humidity in laundry room area, in %
- **T4:** Temperature in office room, in Celsius
- **RH4:**Humidity in office room, in %
- **T5:** Temperature in bathroom, in Celsius
- **RH5:** Humidity in bathroom, in %
- **T6,** Temperature outside the building (north side), in Celsius
- **RH6:** Humidity outside the building (north side), in %
- **T7:** Temperature in ironing room , in Celsius
- **RH7:** Humidity in ironing room, in %
- **T8,** Temperature in teenager room 2, in Celsius
- **RH8:** Humidity in teenager room 2, in %
- **T9:** Temperature in parents room, in Celsius
- **RH9:** Humidity in parents room, in %
- **To,** Temperature outside (from Chievres weather station)
- **Pressure** (from Chievres weather station), in mm Hg
- **RHout,** Humidity outside (from Chievres weather station), in %
- **Wind speed:** (from Chievres weather station), in m/s
- **Visibility:** (from Chievres weather station), in km
- **Tdewpoint:** (from Chievres weather station), Â°C
- **rv1:** Random variable 1, nondimensional
- **rv2:** Random variable 2, nondimensional

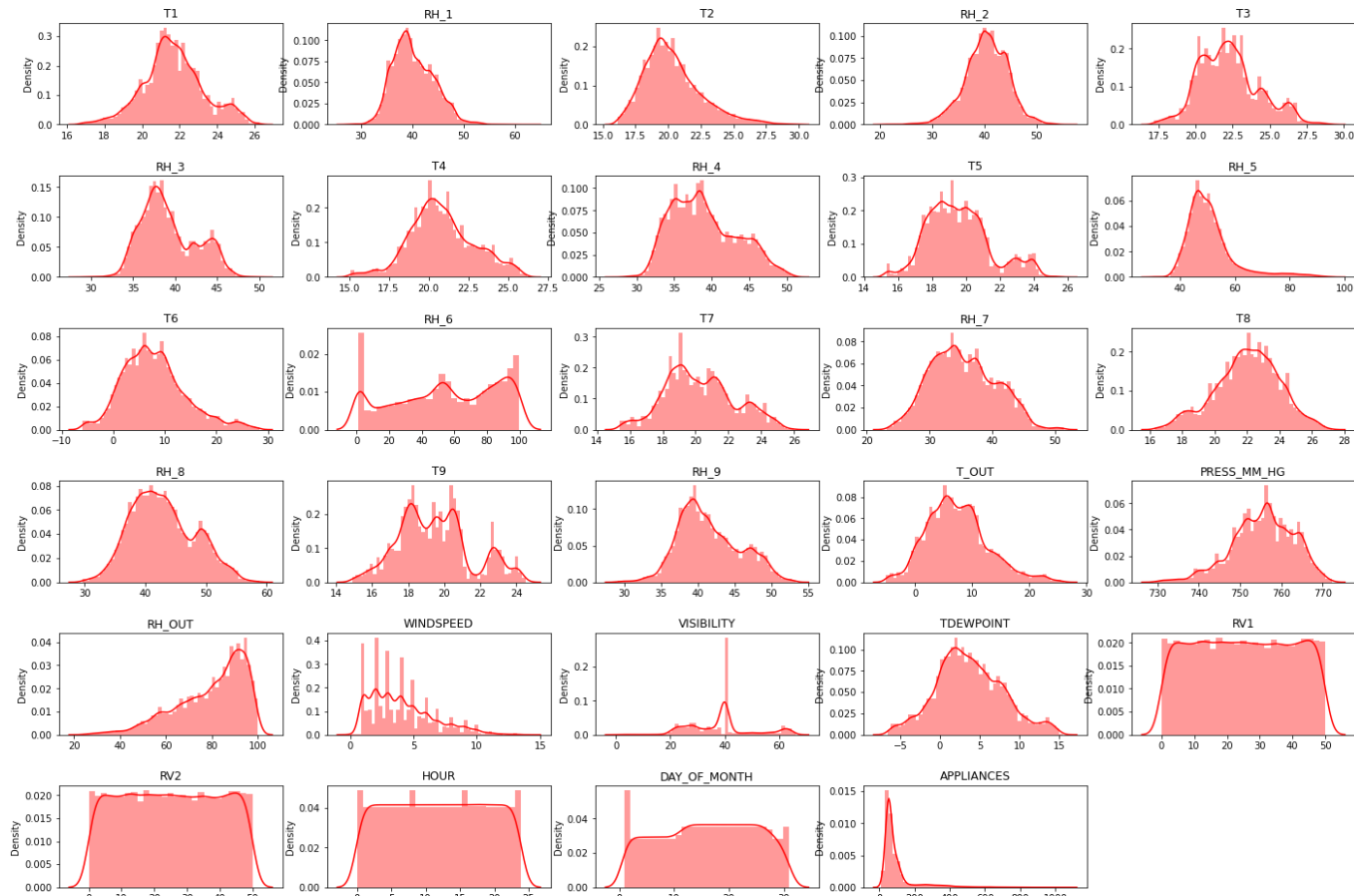
# Dataset

- ❖ Data have 19735 readings and 29 Attributes.
- ❖ No Missing Value
- ❖ No Duplicate Value

Include Categorical Feature as Phase of the Day

Time ( Hours)	Phase of the Day
6 am to 12 pm	Morning
12 pm to 6 pm	Afternoon
6 pm to 12 pm	Evening
12 pm to 6 am	Night

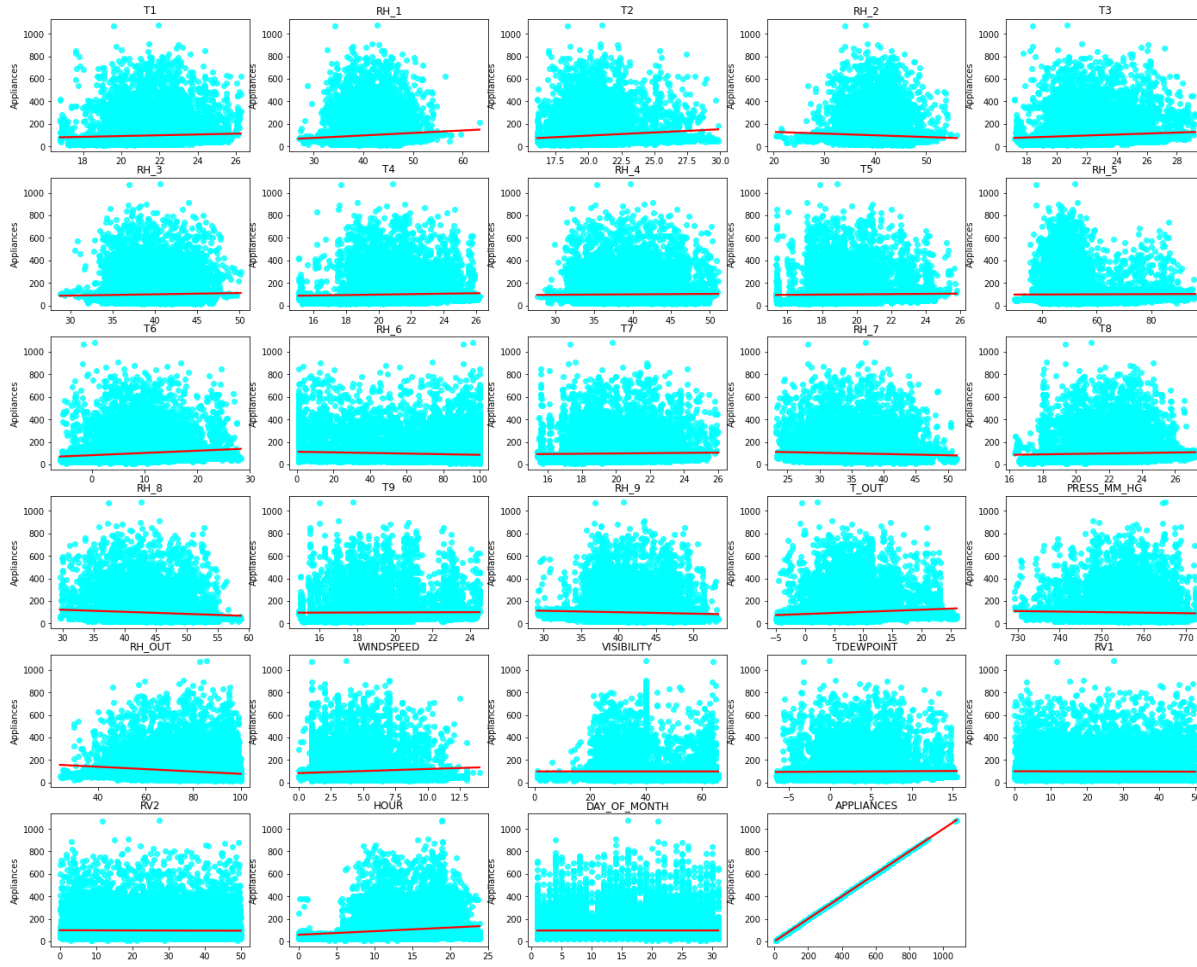
# Exploratory Data Analysis-



1. Temperature and Humidity attributes have a Gaussian-like distribution. The target variable 'Appliances' has a skewed Gaussian distribution indicating a wide range of outliers over the 3rd quartile.

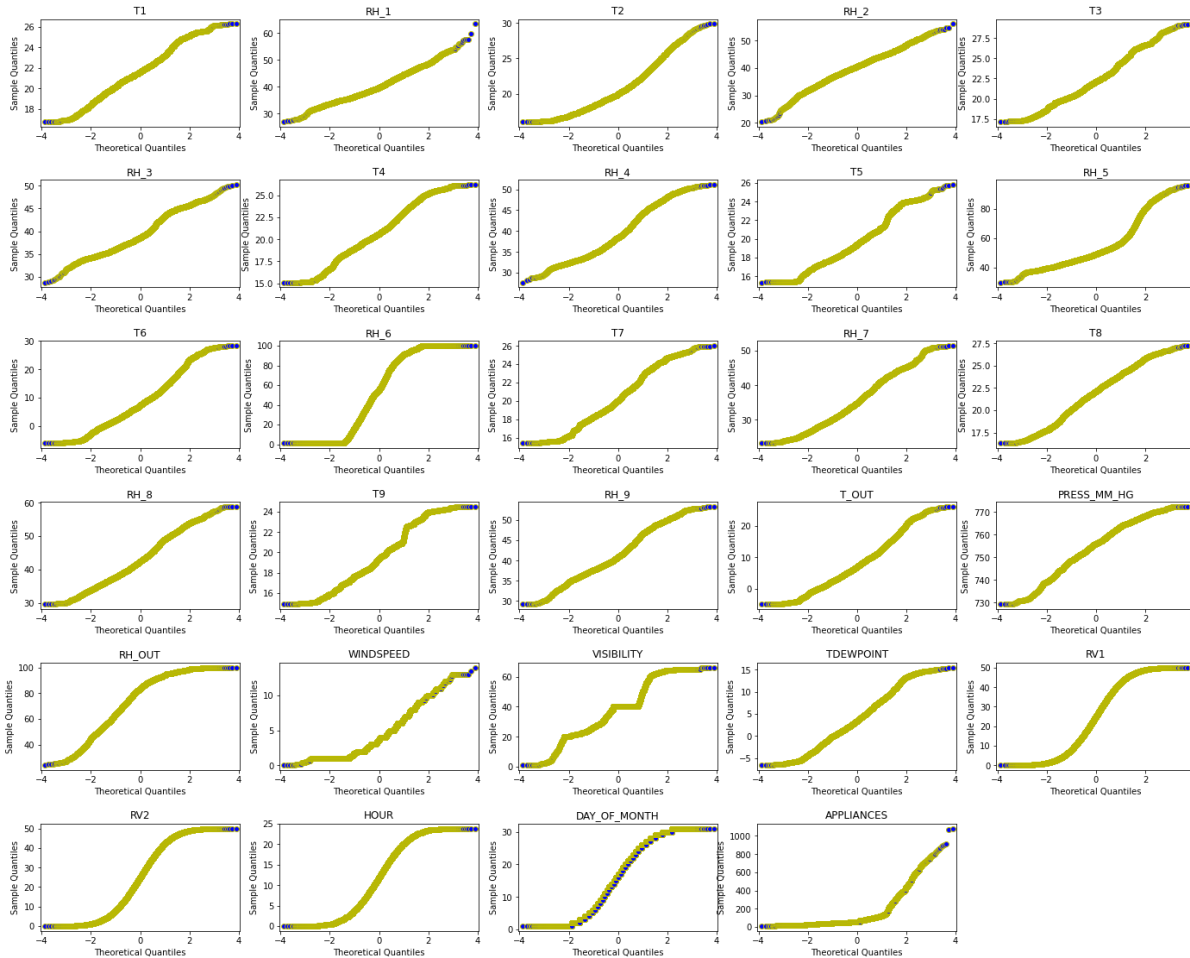
2. Variables hour, rv2, rv1, Appliances, rh6, rv1, T9, wind speed, visibility are not normally distributed. Other variables are seeming to be normally distributed.

# Exploratory Data Analysis-



1. Data is highly nonlinear with dependent variable (Appliances)
2. Regression plot with dependent variables shows very less linear correlation with other variables.

# Exploratory Data Analysis-



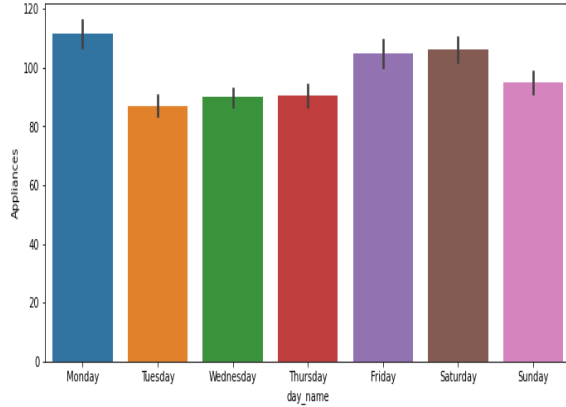
1. The Quantile-Quantile (q-q) plot is a graphical technique for determining if two data sets come from populations with a common distribution.

2. There is large variation in QQ plot of variables hour, rv2, day\_of\_month, appliances, rh6, rv1, T9, T3, windspeed, visibility.

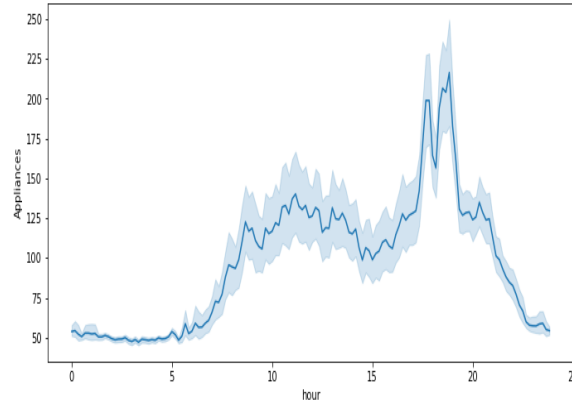


# Exploratory Data Analysis-

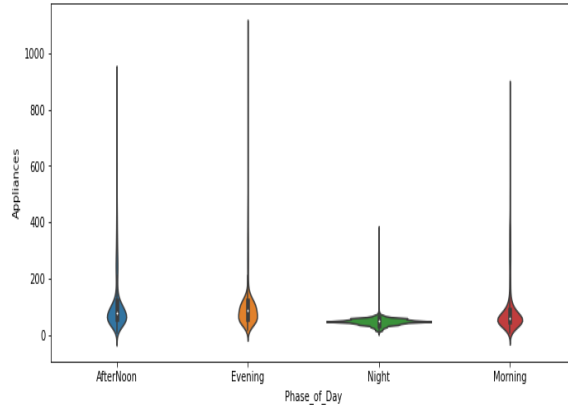
Daywise Energy Consumption



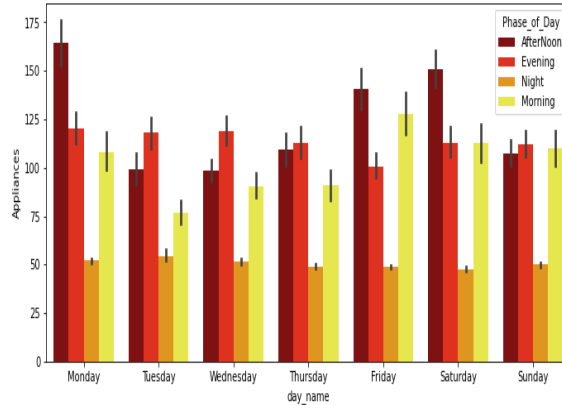
Hourly Day Consumption



Phase of daywise Energy Consumption



Daywise Consumption

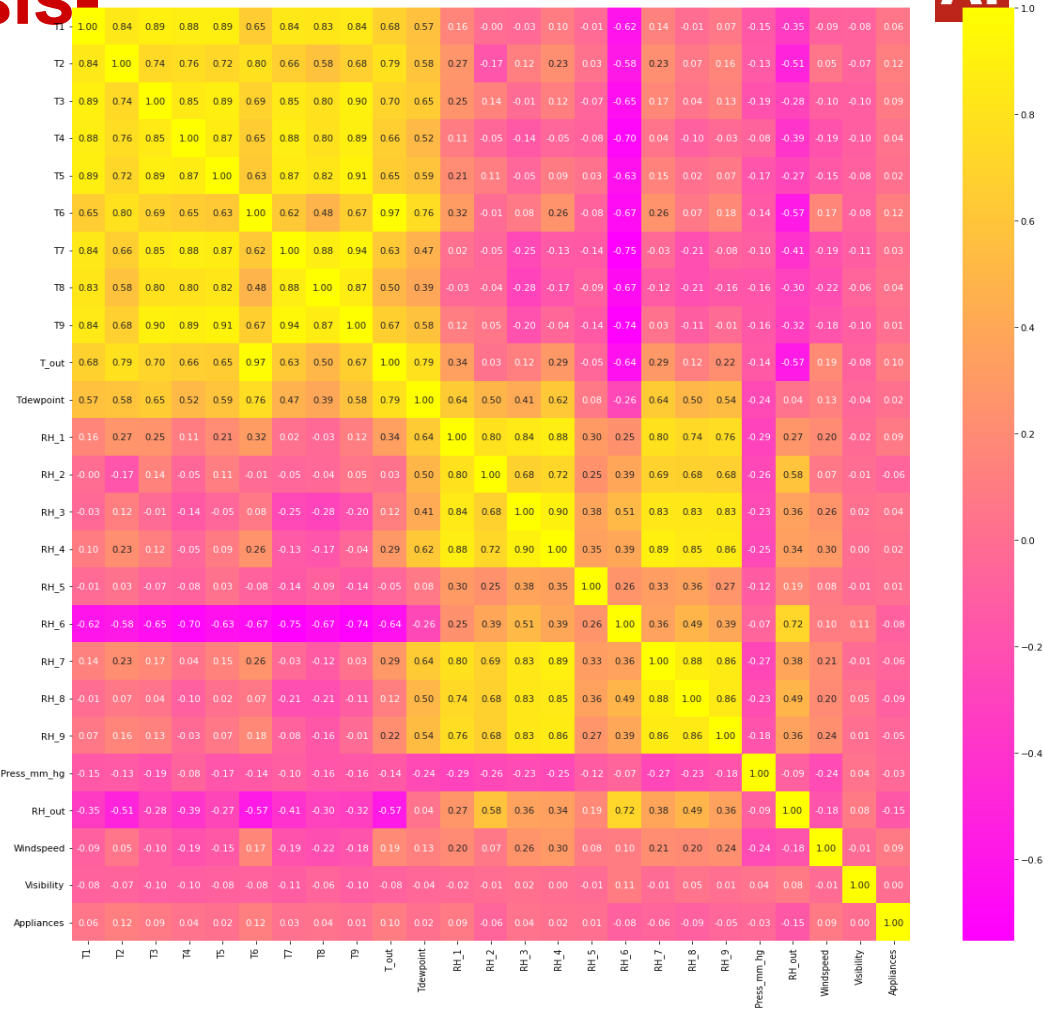


1. Energy consumption at night is less as compared to morning and afternoon on every weekday. Energy consumption is high in the evening.

2. Energy consumption is high on Monday, Friday and Saturday in the afternoon. Most of the appliances consumes around 200 Wh energy

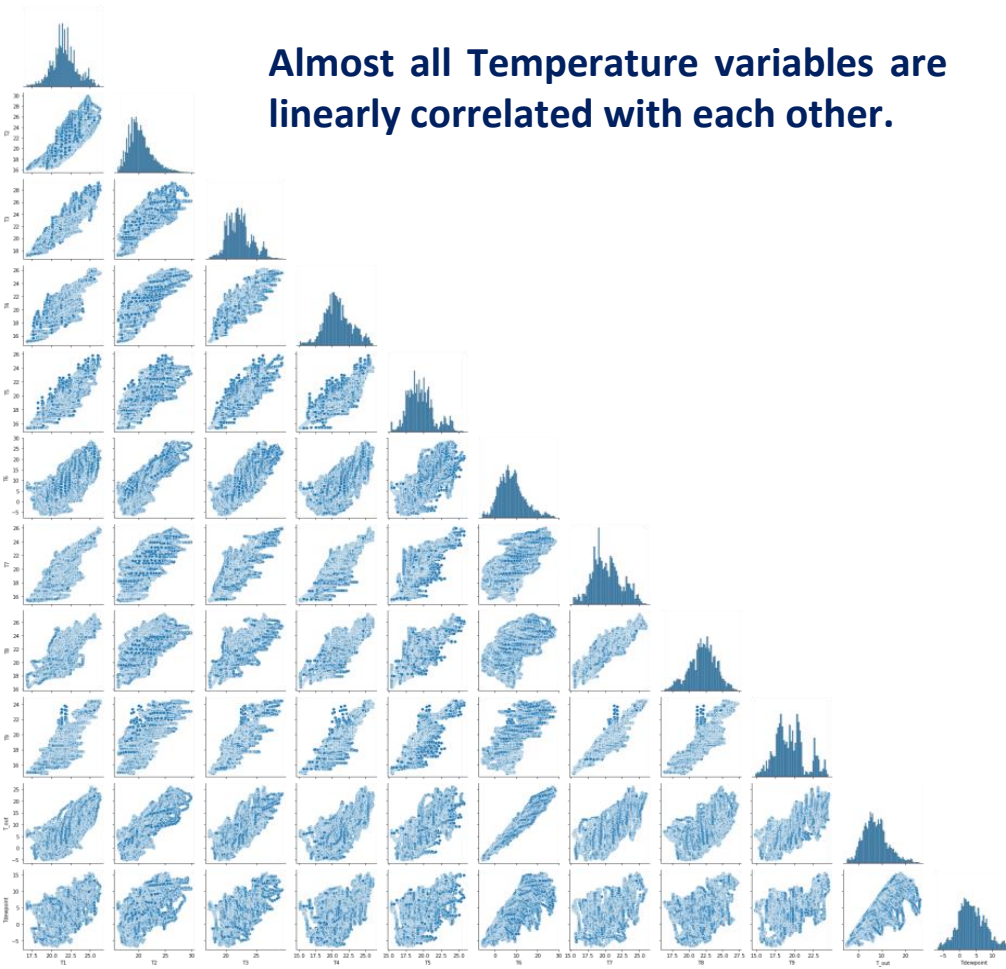
# Exploratory Data Analysis-

1. T9 and T7 are highly correlated as 0.94
2. T\_out and T6 are as 0.97 correlated as 0.97
3. We see strong correlation among temperature variables as change in outside heat can be experienced by all rooms except when changed only by human intervention such as use of thermostat, heaters etc.
4. There is also a strong inverse correlation observed between RH\_6 and all temperature features. This is because RH\_6 is the outside humidity.
5. As air temperature increases, air can hold more water molecules, and its relative humidity decreases

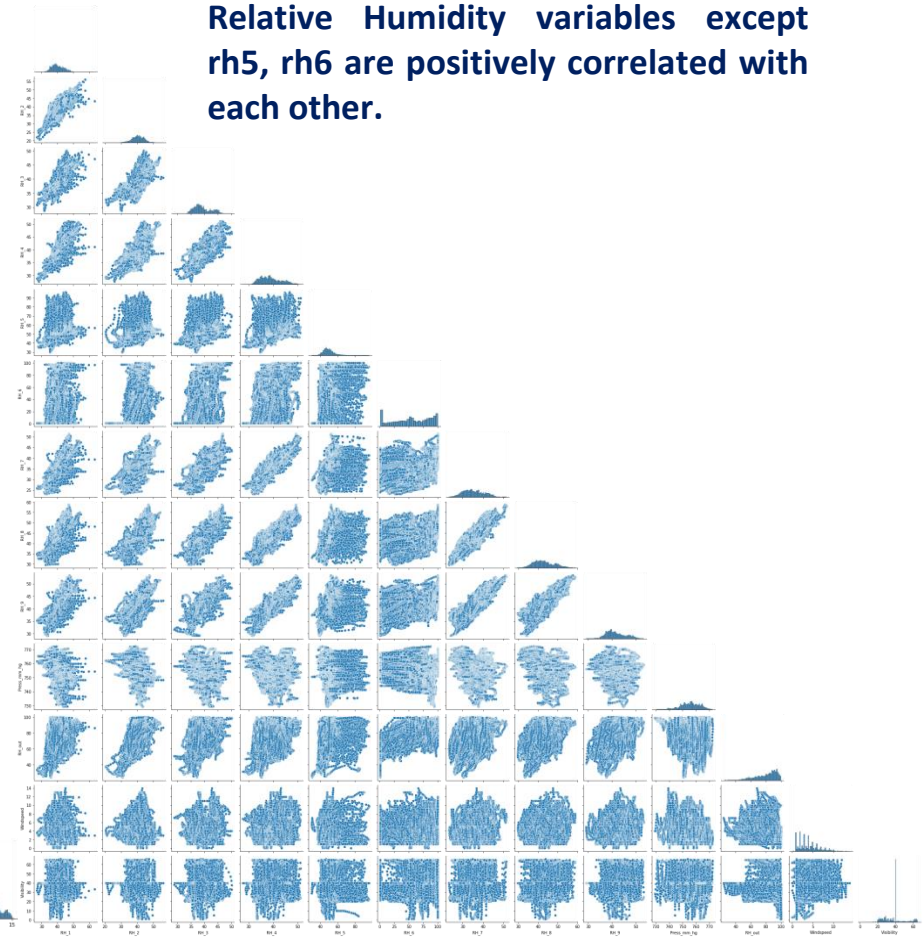


# Exploratory Data Analysis-

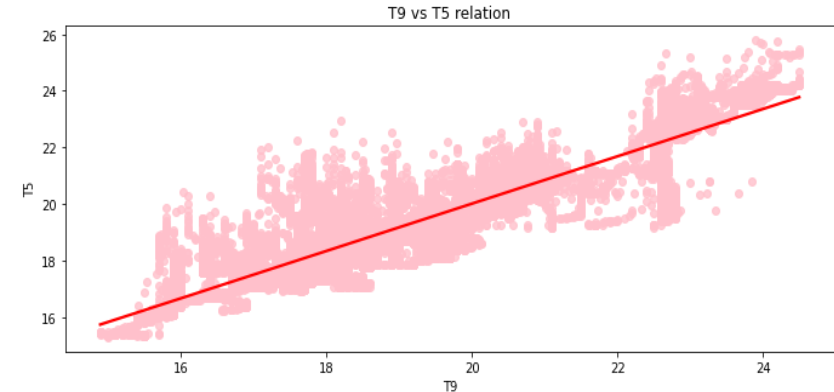
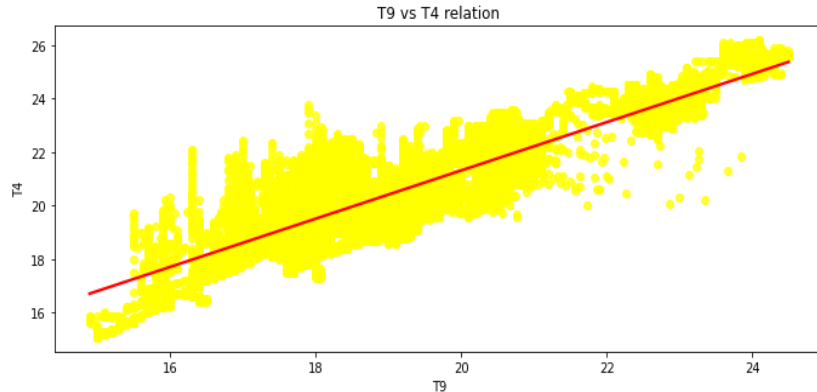
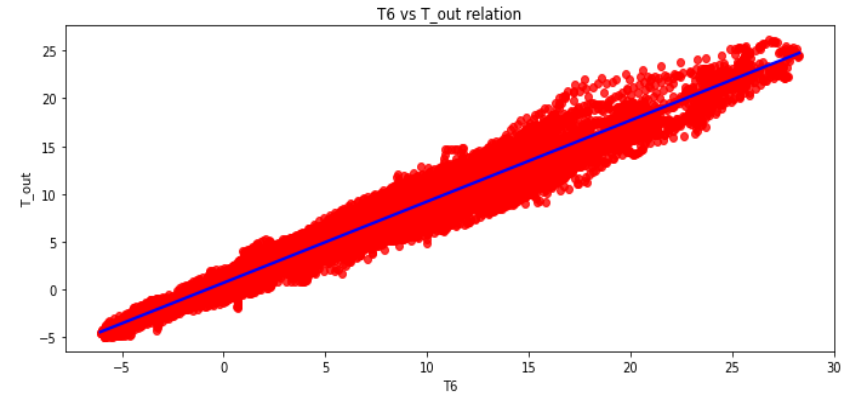
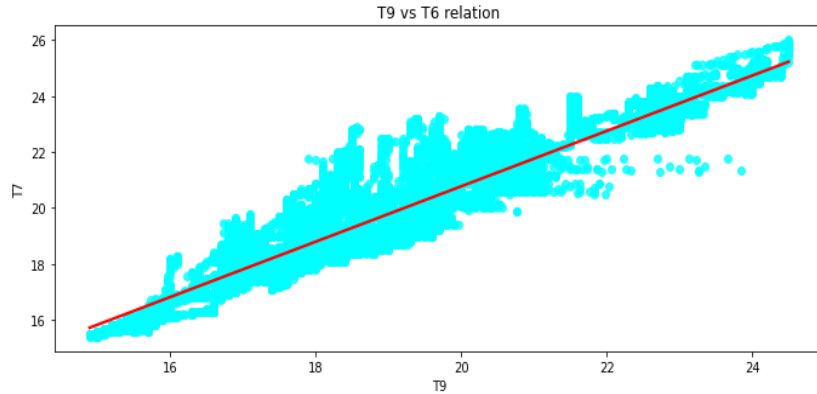
Almost all Temperature variables are linearly correlated with each other.



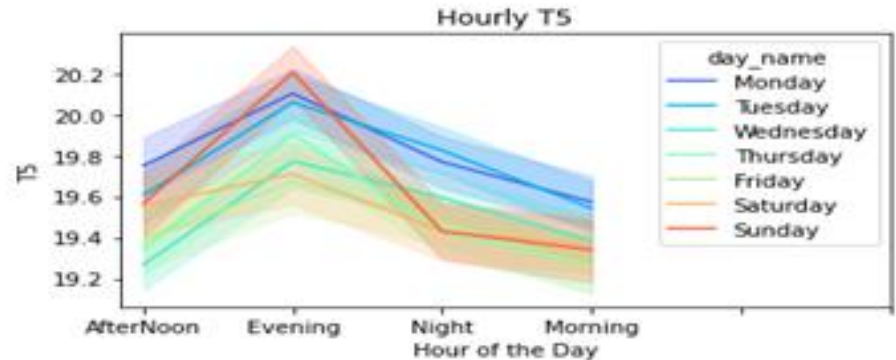
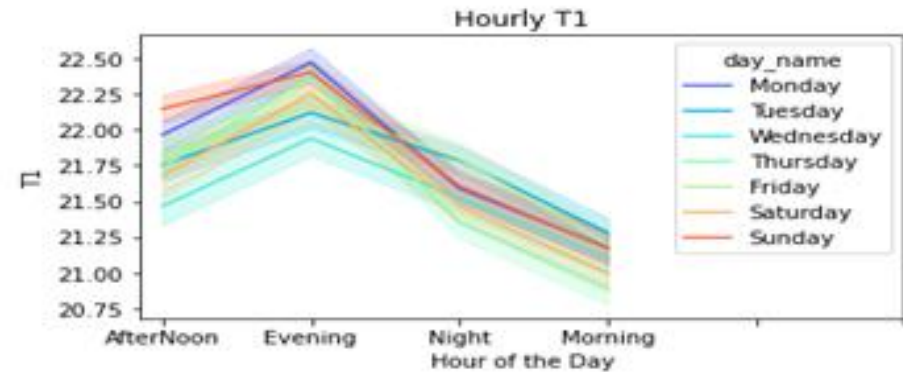
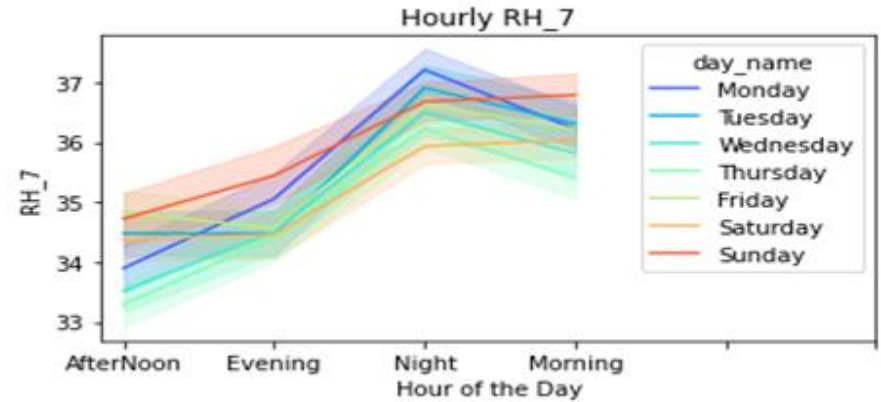
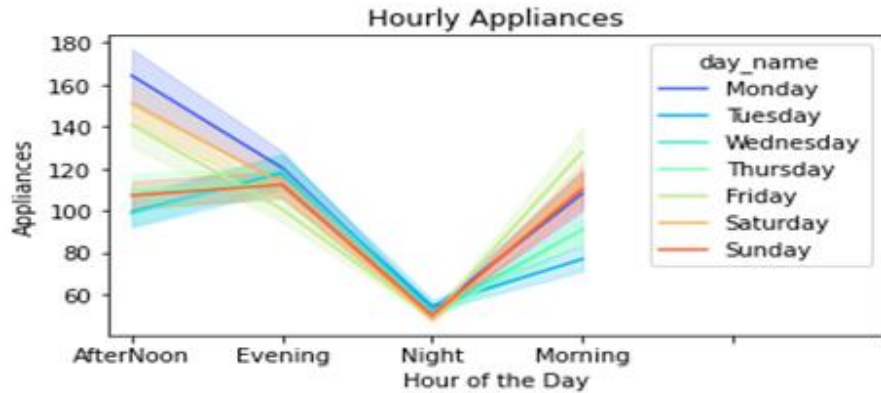
Relative Humidity variables except rh5, rh6 are positively correlated with each other.



# Exploratory Data Analysis-

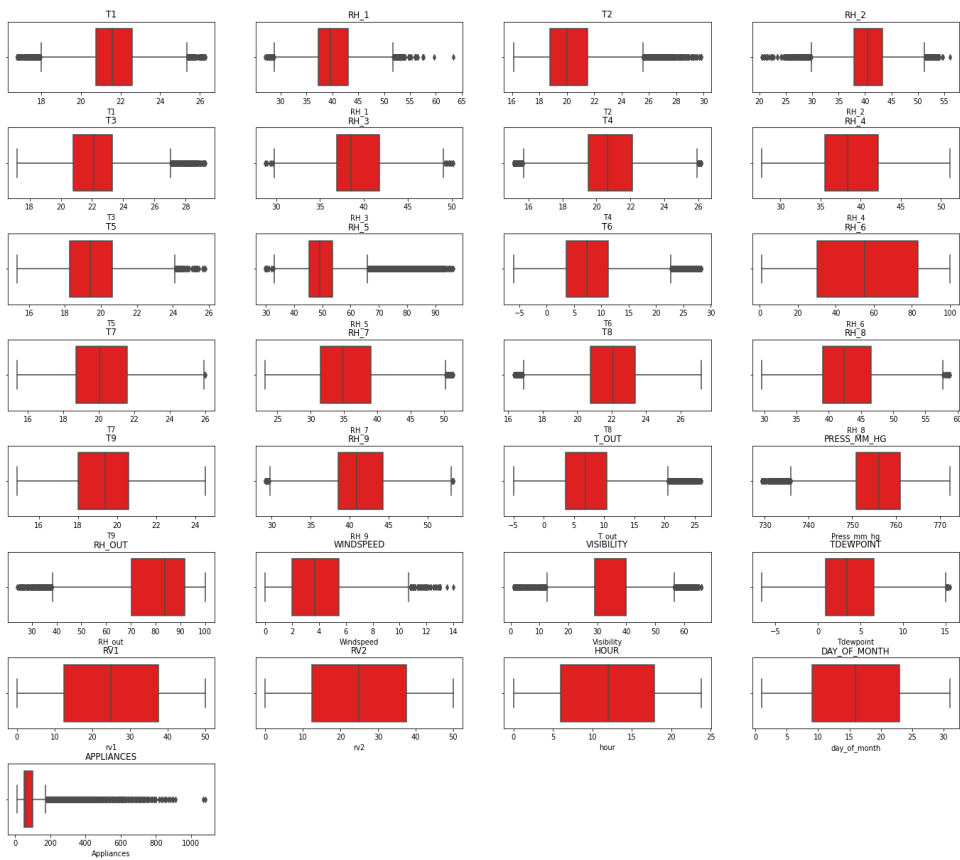


# Exploratory Data Analysis-

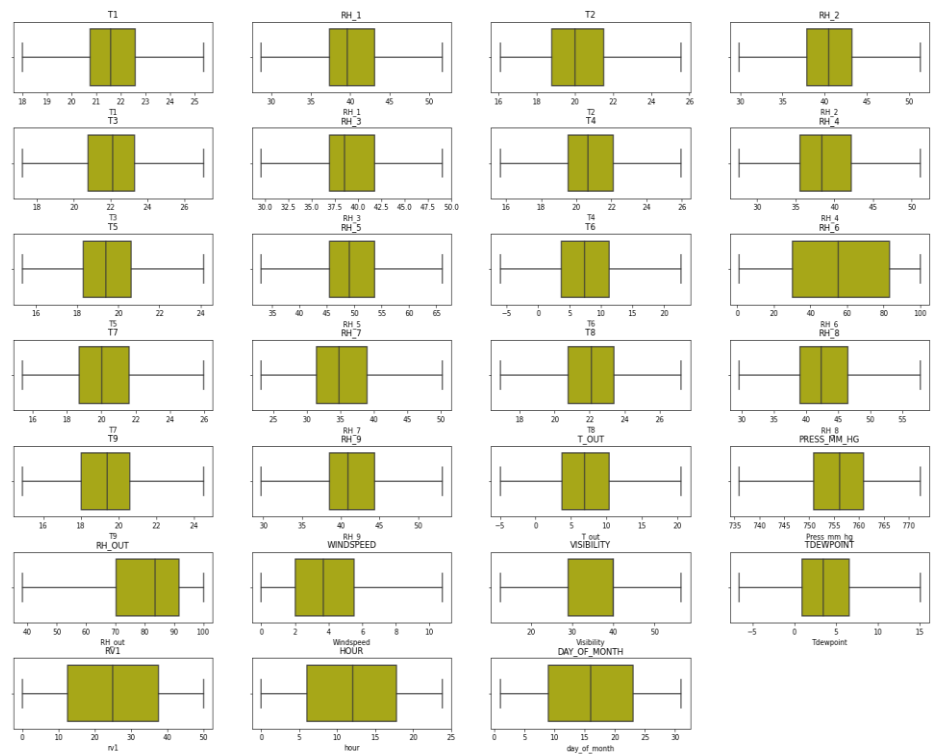


# Feature Engineering and Selection

## Before removing Outliers



## After removing Outliers



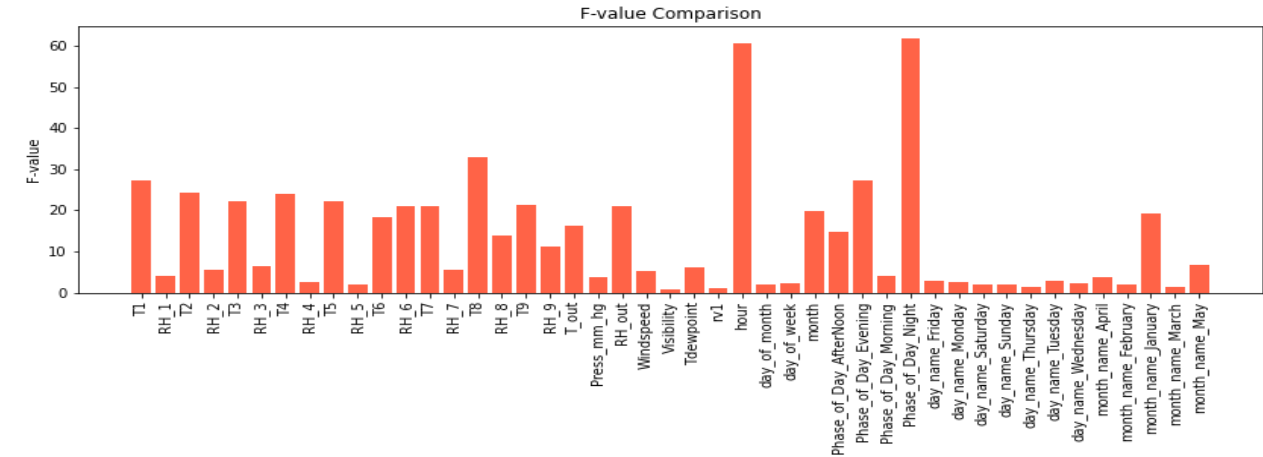
# Feature Selection

## ANOVA F Value

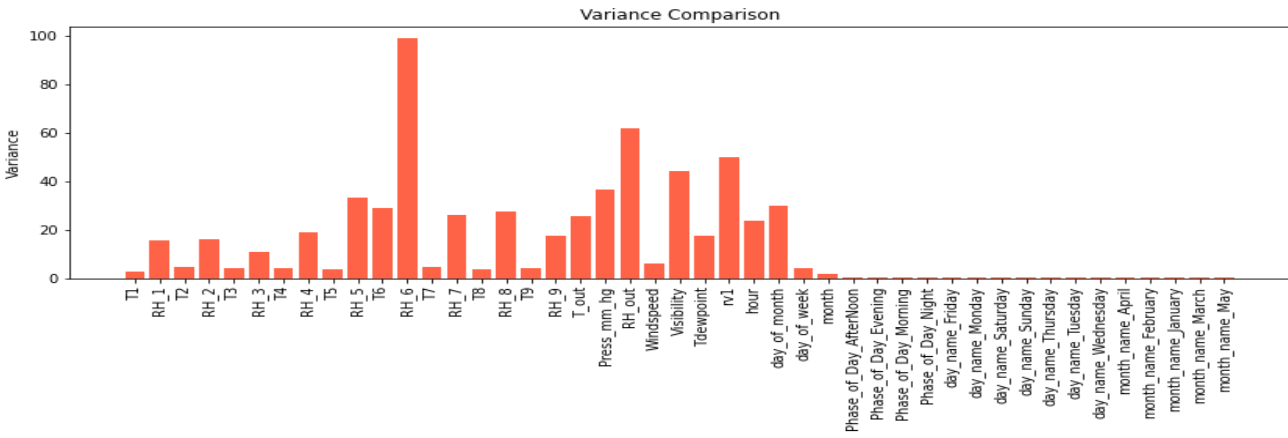


### Check Multicollinearity

	variables	VIF
25	rv1	1.002546
23	Visibility	1.050790
0	Appliances	1.213415
10	RH_5	1.534392
20	Press_mm_hg	1.563047
27	day_of_month	1.711111
22	Windspeed	1.717932
18	RH_9	7.886518
7	T4	9.923554



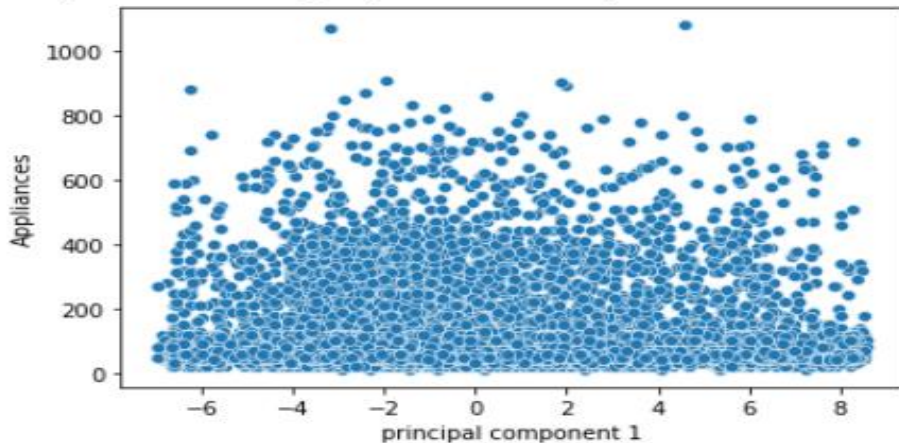
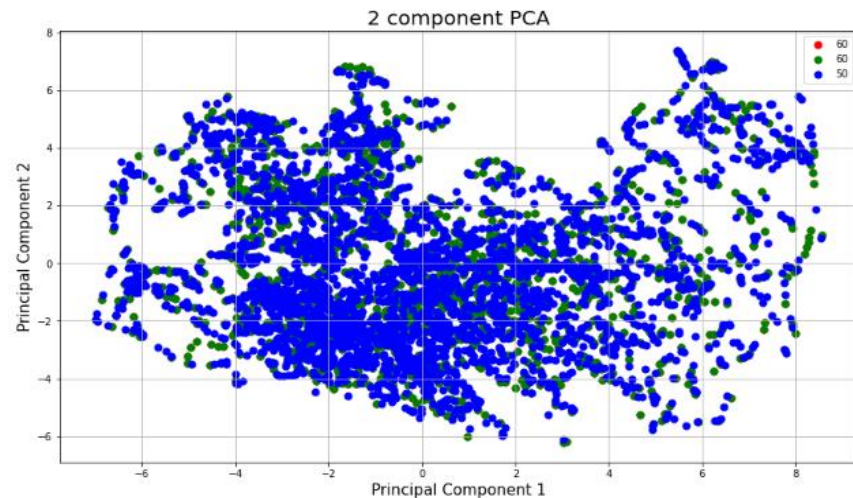
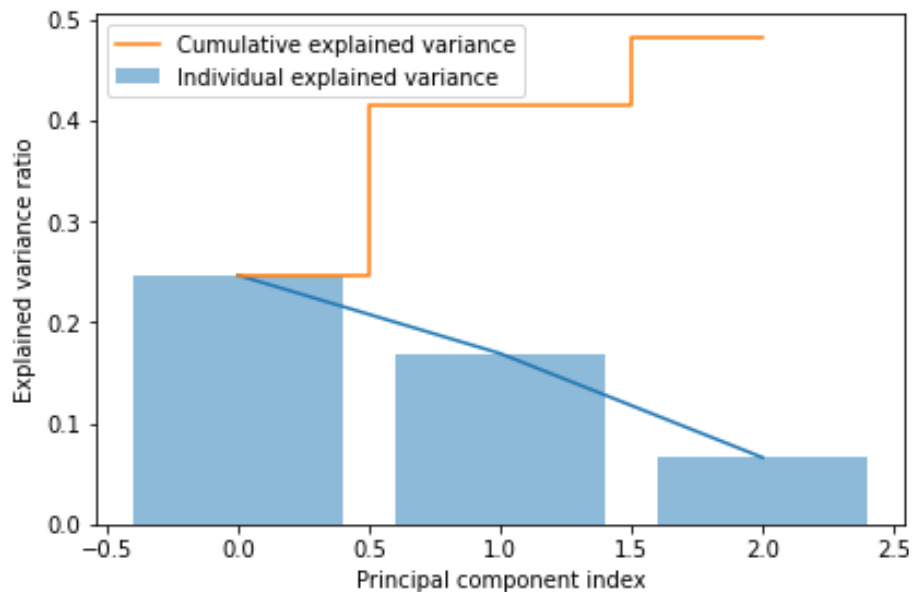
### Variance Threshold





# Feature Selection

## Principal Component Analysis





## AI

**Training Data 80%**  
(Rows, Columns)

**Test data 20%**  
(Rows, Columns)

[illegible]

# Developing Machine Learning Model

## Evaluation Metrics

Divide by total Number of Data Points

Actual Output

Predicted Output

$$MAE = \frac{1}{N} \sum |y - \hat{y}|$$

Sum Of

Absolute Value of residual

$$MSE = \frac{1}{n} \sum \underbrace{(y - \hat{y})^2}_{\text{The square of the difference between actual and predicted}}$$

**RMSE =  $\sqrt{MSE}$**

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

$$R^2 \text{ Squared} = 1 - \frac{SSr}{SSm}$$

SSr = Squared sum error of regression line

SSm = Squared sum error of mean line

$$R_a^2 = 1 - \left[ \left( \frac{n-1}{n-k-1} \right) \times (1 - R^2) \right]$$

where:

n = number of observations

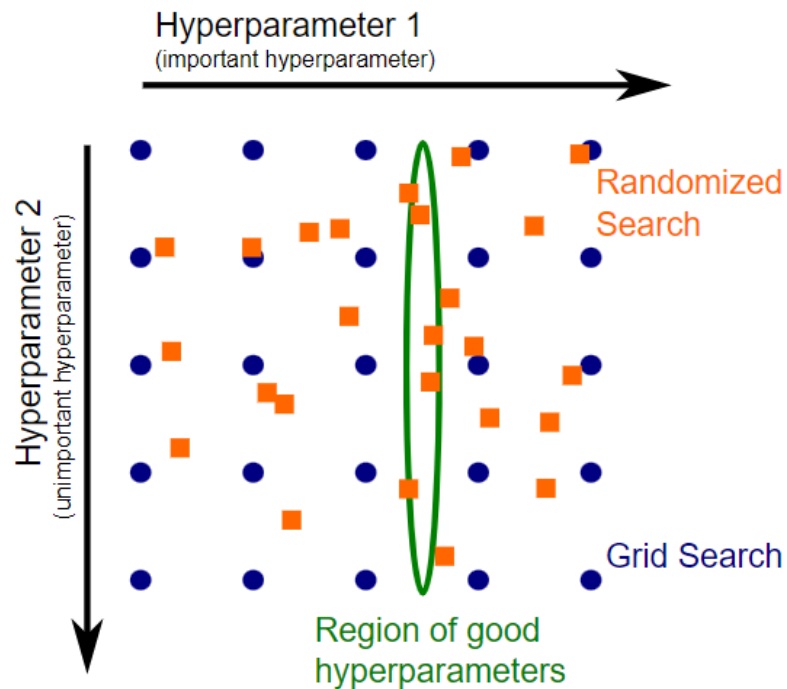
k = number of independent variables

$R_a^2$  = adjusted  $R^2$

# Cross Validation, Hyper Parameter Tuning ( GridSearchCV, RandomizedSearchCV)

Fitting 5 folds for each of 10 candidates, totalling 50 fits

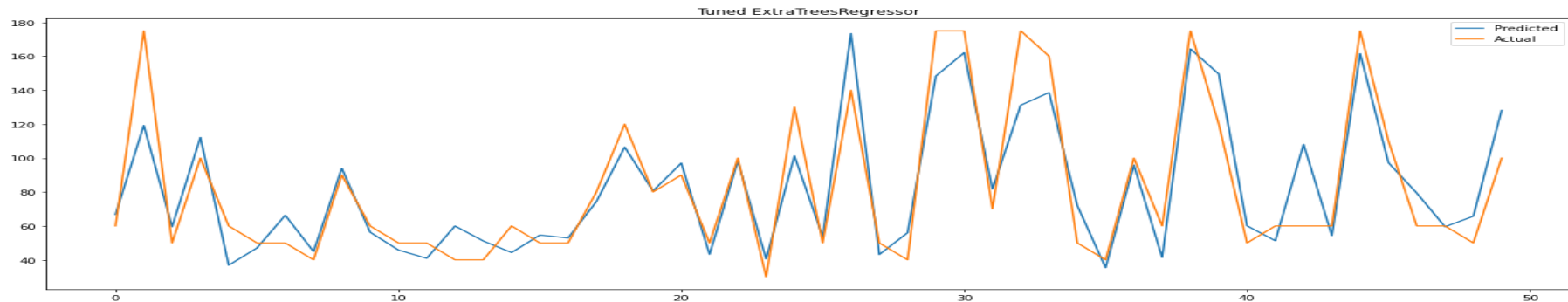
```
RandomizedSearchCV(estimator=ExtraTreesRegressor(random_state=5), n_jobs=-1,
                   param_distributions={'bootstrap': [True, False],
                                       'criterion': ['squared_error'],
                                       'max_depth': [80, 100, None],
                                       'max_features': ['log2', 'sqrt'],
                                       'n_estimators': [10, 1000, 100]},
                   scoring='r2', verbose=10)
```



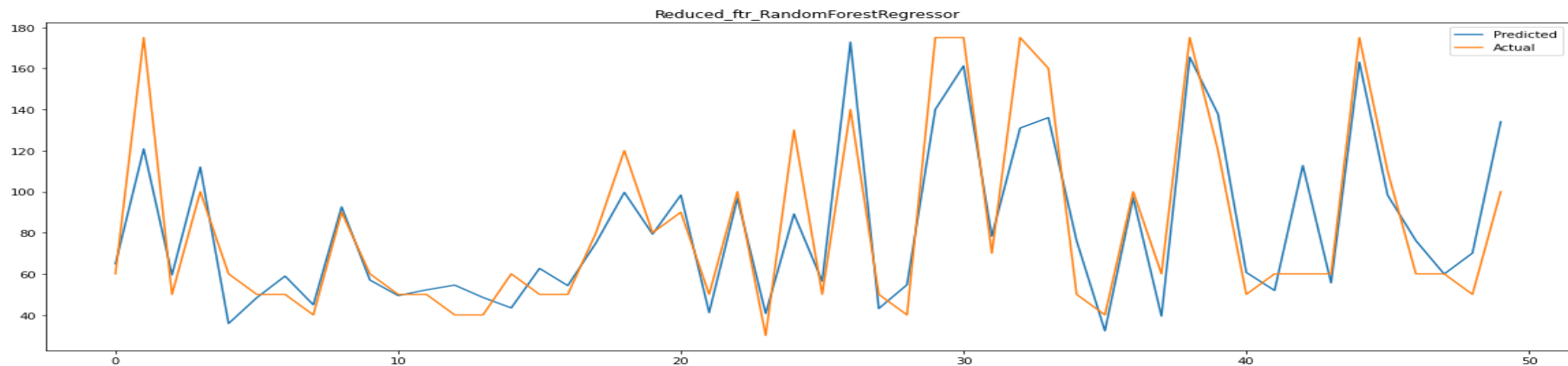
```
({'n_estimators': 1000,
  'max_features': 'sqrt',
  'max_depth': 100,
  'criterion': 'squared_error',
  'bootstrap': False},
 0.7455214226884088)
```

# Selecting Optimal Model

## Extra Tree Regressor with Hyperparameter Tuning



## Extra Tree Regressor with Feature reduction and Hyperparameter Tuning

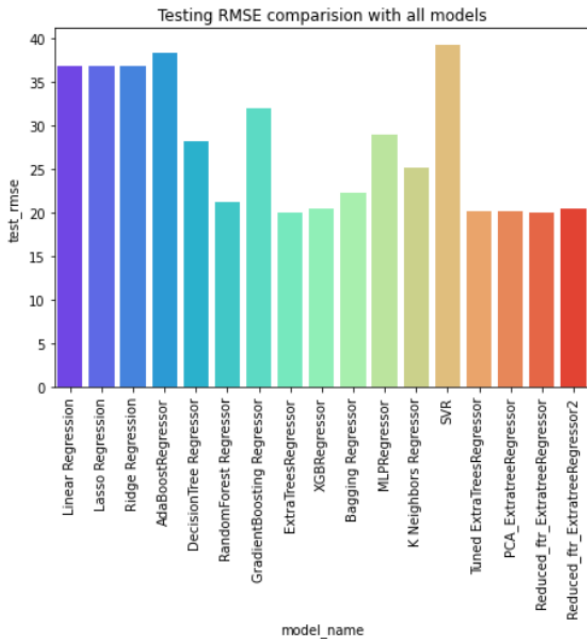
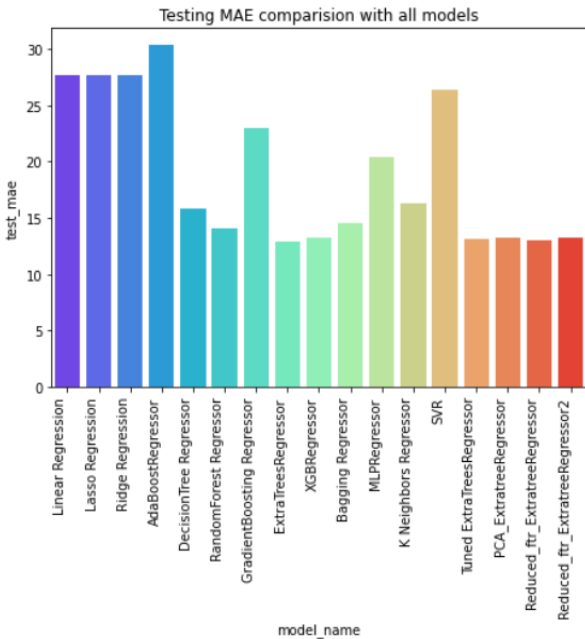
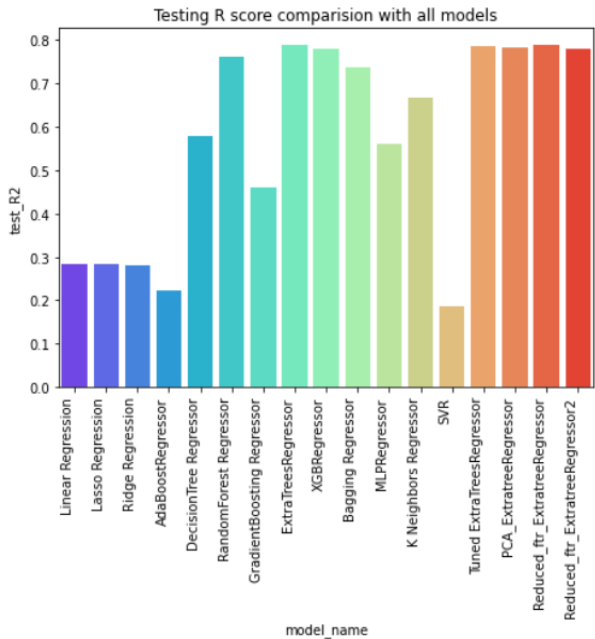


AI

# Model Performance Report

	model_type	model_name	test_mse	test_rmse	test_mae	test_R2	test_adjusted R2	training_R2
0	Linear	Linear Regression	1360.106641	36.879624	27.660606	0.282230	0.274693	0.303539
1	Linear	Lasso Regression	1360.115140	36.879739	27.660610	0.282225	0.274689	0.303539
2	Linear	Ridge Regression	1360.833017	36.889470	27.664881	0.281846	0.274306	0.303503
3	Ensemble Method	AdaBoostRegressor	1475.788565	38.415994	30.317404	0.221181	0.213004	0.228933
4	Tree	DecisionTree Regressor	799.461616	28.274752	15.838612	0.578099	0.573670	1.000000
5	Ensemble Method	RandomForest Regressor	449.257687	21.195700	14.050621	0.762913	0.760424	0.965790
6	Ensemble Method	GradientBoosting Regressor	1024.710152	32.011094	22.952282	0.459229	0.453551	0.503818
7	Ensemble Method	ExtraTreesRegressor	400.521013	20.013021	12.911211	0.788633	0.786413	1.000000
8	Ensemble Method	XGBRegressor	416.246358	20.402116	13.215798	0.780334	0.778028	1.000000
9	Ensemble Method	Bagging Regressor	495.999557	22.271047	14.531670	0.738246	0.735497	0.951160
10	NN Method	MLPRegressor	834.791996	28.892767	20.396612	0.559454	0.554829	0.603148
11	Neighbours	K Neighbors Regressor	629.741069	25.094642	16.322777	0.667666	0.664177	0.800531
12	SVM	SVR	1540.013910	39.243011	26.354987	0.187287	0.178754	0.221656
13	Ensemble Method	Tuned ExtraTreesRegressor	405.826868	20.145145	13.105933	0.785833	0.783529	0.999999
14	Ensemble Method	PCA_ExtiratreeRegressor	408.505201	20.211512	13.192513	0.784419	0.781932	1.000000
15	Ensemble Method	Reduced_ftr_ExtiratreeRegressor	402.178913	20.054399	13.057385	0.787758	0.786241	1.000000
16	Ensemble Method	Reduced_ftr_ExtiratreeRegressor2	416.254120	20.402307	13.283000	0.780330	0.779267	1.000000

# Model Performance Report



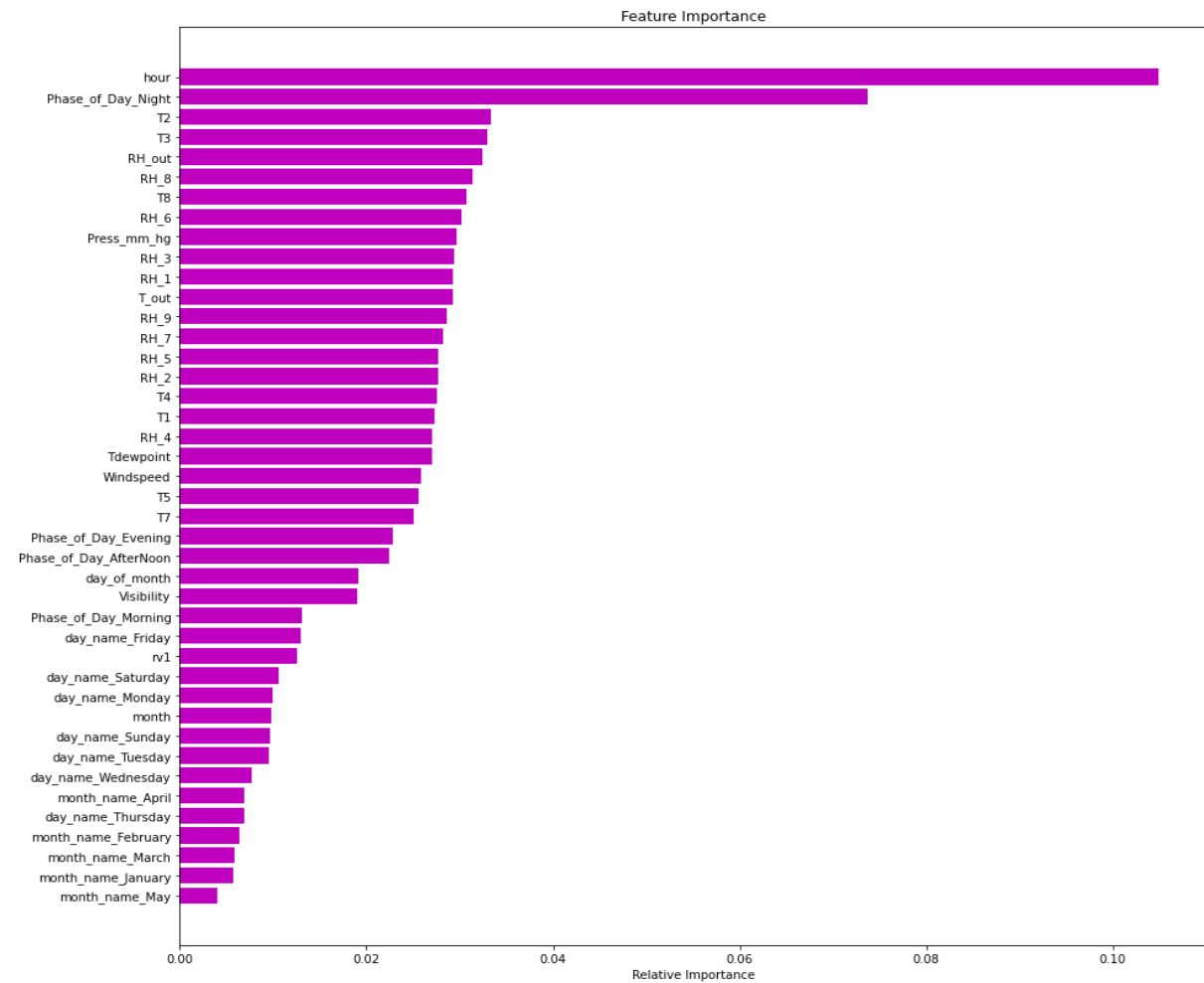
# Model Explainability

- Model Explainability refers to the concept of being able to understand the machine learning model Importance: Feature Importance refers to techniques that assign a score to input features based on how useful they are at predicting a target variable. Feature importance scores provide insight into the data and the deployed model.

## Feature Importance

By looking at the Feature Importance graphs and the contribution chart from ELI5, we can gather that the appliance energy consumption largely depends on the 'T1','RH\_1','T2','RH\_2','T3','RH\_3','T4','RH\_4','T5','RH\_5','RH\_6','T7','RH\_7','T8','RH\_8','RH\_9','T\_out','Press\_mm\_hg','RH\_out','Windspeed','Tdewpoint',hour

# Feature Importance





# LIME (Local Interpretable Model-Agnostic Explanations) and Eli5

- LIME, the acronym for Local Interpretable Model-Agnostic Explanations, is a technique that approximates any black box machine learning model with a local, interpretable model to explain each individual prediction.

y (score **48.600**) top features

Contribution?	Feature	Value
+78.651	<BIAS>	1.000
+1.151	T2	0.235
+0.669	rv1	1.527
+0.669	T1	0.012
+0.438	RH_5	0.108
+0.429	RH_4	0.912
+0.210	T3	-0.879
+0.209	T5	-0.260
+0.106	T4	-0.627
+0.104	day_of_month	-0.005
+0.094	RH_7	1.235
+0.092	Tdewpoint	-0.897
+0.075	day_name_Thursday	-0.413
+0.059	day_name_Monday	-0.404
+0.034	day_name_Saturday	2.514
+0.032	RH_2	-0.236
+0.024	day_name_Friday	-0.413
-0.024	RH_3	0.815
-0.030	month_name_March	-0.544
-0.040	month_name_May	-0.491
-0.063	day_name_Wednesday	-0.412
-0.077	month_name_February	-0.518
-0.168	month_name_April	-0.526
-0.188	day_name_Tuesday	-0.416

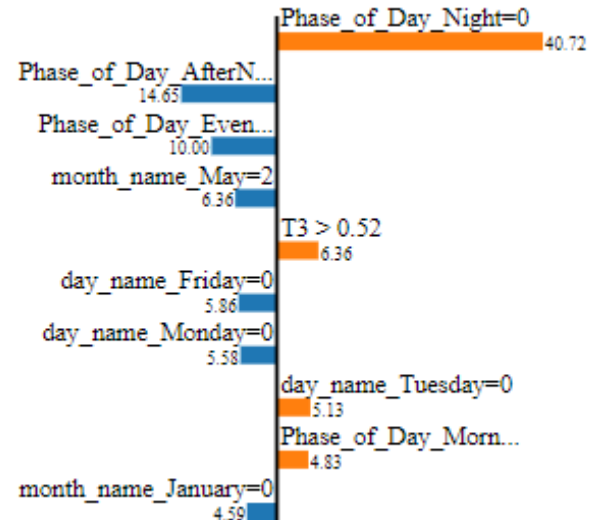
Predicted value

56.45 (min) 61.10 200.43 (max)

Actual Value 60

negative

positive



# Limitations

1. One of the main limitations of this study is that the analysis was done for only one house.
2. Important information could be found when analyzing several houses, and other relationships can be studied with appliances' energy consumption in combination with: occupant's age, number of occupants, ownership of pets, building's geometry etc.
3. Another research limitation is the length of continuous analyzed data. Different energy use patterns can potentially be found depending on the season of the year.
4. Regarding the weather station, the predictions of appliances energy use could probably be better if the weather station was closer to the house.
5. This research has not looked into the problem of optimal location of the wireless indoor sensors for improvement of the energy prediction. It is also possible that more sensors and better sensor accuracy could help to improve the energy prediction.

# Conclusion

1. The household appliance energy consumption prediction models based on Linear Regression, Lasso Regression, Ridge Regression, MLP Regressor, DecisionTree Regressor Random Forest Regressor, Adaptive Boosting Regressor, Gradient Boosting Regressor, Bagging Regressor, XGBoost, K Neighbors Regressor and Linear SVM are explored.
2. The variables T6 and T\_out , T9 and T7 has high correlation with each other hence we have dropped T6 and T9 to avoid Multicollinearity. Upon appropriate pre-processing and fitting the models, we compare and evaluate the best model with lowest error and the highest R-squared score.
3. When evaluating the influence of Random Variable attribute the linear models have assigned near zero weights to the random variable, negating its influence in prediction of the target variable.
4. Extra Tree Regressor was found to be the best performing model with an R-squared score of 0.7877.
5. After optimizing the hyperparameters of the Extra Tree Regressor, doing principle Component Analysis, its R-squared score increased from 0.7833 to 0.7877.

# Conclusion

- 6. We find that this model's predictions are mainly contributed by the hour, 'T1', 'RH\_1', 'T2', 'RH\_2', 'T3', 'RH\_3', 'T4', 'RH\_4', 'T5', 'RH\_5', 'RH\_6', 'T7', 'RH\_7', 'T8', 'RH\_8', 'RH\_9', 'T\_out', 'Press\_mm\_hg', 'RH\_out', 'Windspeed', 'Tdewpoint'. Temperature and Relative Humidity of kitchen, living room, laundry room, Ironing room, outside surrounding are playing important role in Energy Prediction.
- 7. Data from a wireless sensor network that measures humidity and temperature has been proven to increase the prediction accuracy. The data analysis showed that data from the kitchen, laundry room, living room and bathrooms had the most important contributions. Data from the other rooms also helps in the prediction. When looking at the appliances in each room, it can be seen that the laundry, kitchen and living rooms would be expected to have the highest contributions because of the equipment present. The prediction of appliances' consumption with data from the wireless network indicates that it can help to locate where in building the main appliances' energy consumption contributions are found.

# Conclusion

- 8. When using all the predictors the light consumption was ranked highly. However, when studying different predictor subsets, removing the light consumption appeared not to have a significant impact. This may be an indication that other features are correlated well with the light energy consumption.
- The possible explanation for why the pressure has a strong prediction power may be related to its influence on the wind speed and higher rainfall probability which could potentially increase the occupancy of the house.
- 9. As this dataset has a time component to it, we believe that better performances can be achieved by using Time Series Analysis concepts.

## Future Scope/ Suggestions

- This study has found curious relationships between variables. Future work could include considering weather data such as solar radiation and precipitation. Also, occupancy and occupant's activity information could be useful to improve the prediction and find its relationship with other parameters (exterior weather for example). The wireless sensors could also measure CO<sub>2</sub> and noise to help in the prediction and to track the occupant's movement from room to room and time spent in each room.

## References:

1. Luis M. Candanedo, Véronique Feldheim, Dominique Deramaix, 'Data driven prediction models of energy use of appliances in a low-energy house', Energy and Buildings, Thermal Engineering and Combustion Laboratory, University of Mons, Rue de l' Epargne 56, 7000 Mons, Belgium.
2. <https://www.geeksforgeeks.org/python-programming-language/>
3. [https://scikit-learn.org/stable/supervised\\_learning.html#supervised-learning](https://scikit-learn.org/stable/supervised_learning.html#supervised-learning)
4. <https://scikit-learn.org/stable/modules/ensemble.html#forest>

THANK YOU