# ASC61013 Data Modelling and Machine Intelligence

# Course Work – Solar Power Consumption and Generation

Ankur Singh Gulia

Registration Number: - 220181437

# *Content*

# 1. Introduction to Coursework

The data set provided consists of 50,392 instances with 29 features and 1 meta – attribute in the raw CSV data file named CW_dataset.

The features in the raw CSV datasets are related to the effects on solar power generation and consumption of the energy generated. The meta-attribute present in the raw dataset is time which is spread over 7 days starting from 1st Jan,2016 to 7th Jan 2016. It depicts the energy generation and consumption varying over 7 days with dependency on several weather features like temperature, dewpoint, wind bearing, etc.

This coursework aims to analyse the data and predict the future usage of energy depending on the given data using various regression models. With the understanding of the domain, define a correlation between the features with the target feature. With the help of a PCA or hierarchical clustering analysis predicting which weather feature to keep within the study. Give out a machine learning model to predict the value of energy used from at least two features.

Also, the use of feature engineering to classify the usage of energy based on three classes: Low, Medium, and High. Also, showing which appliance uses most energy in the above categories.

Once the predictions are confirmed, cross-validation of the machine learning pipeline is supposed to be the next step. The application of learning curves and classification evaluation metrics, prove that the pipeline is effective at preventing underfitting as well as overfitting.

Once the learning curve is complete, the understanding of the mathematical peculiarities for different models used can be described along with the reason which regression and classification is best suited for this machine learning model.

# 2. Domain Analysis – Energy Consumption and Generation

Solar energy is the energy that is generated by the sun through the process of nuclear fusion. It is a renewable and clean source of energy that can be used to generate electricity or provide heat. Solar panels, which are made up of photovoltaic cells, are used to capture the sun's energy and convert it into usable electricity. Solar energy has the potential to provide a significant portion of the world's energy needs, and its use is growing rapidly as technology advances and costs decrease.

It has been more common to use solar energy now. There have been quite a few prominent factors that affect the generation for the energy. These factors that affect the generation include:

- Temperature - Measure of the thermal energy of a system or substance.
- Apparent Temperature - Measure of how hot or cold it feels to the human body based on the combination of the air temperature, relative humidity, and wind speed.
- Wind Bearing - also known as wind direction or wind direction angle, is the direction from which the wind is blowing.

- Wind Speed – A measure of how fast the wind is blowing.
- Humidity - Humidity is a measure of the amount of water vapor in the air.
- Precipitation intensity - Measure of the rate at which precipitation (e.g., rain, snow, sleet, or hail) is falling at a given moment.
- Precipitation probability - Measure of the likelihood that precipitation will occur at a given location over a specified period of time.
- Visibility - The distance that one can see clearly.
- Pressure - Measure of the force applied to a given area.
- Dew point - The dew point is the temperature at which moisture in the air will condense into liquid.

The above-mentioned factors have been used as the features to predict the model. These factors affect the generation of energy. As for the consumption, it can be relative to various things but for a house, usage can be categorised easily like furnace, kitchen, barn, well, garage, living room, etc. These categories have been taken as the features for consumption of energy which sum up to be 15 features.

According to the data, the weather features that affect the most can be figured out after the data pre-processing is complete with the help of correlation. At present the data is given out as :



*Figure 1 Data Screening*

From the quick pre analysis and understanding of the domain, it can be stated that the weather features affecting the generation of the energy would be Pressure and temperature whereas for the usage, the features to be considered are Pressure, dewpoint and apparent temperature.

But this is only initial layout of the analysis based on the domain analysis. It may be very well that some of the features are cleared to be not of use in the in depth analysis.

## 3. Data Pre-processing

The dataset generally has certain unwanted or non-interactive data which does not allow the pipeline to work in accordance with the model which is required to work the dataset and features.
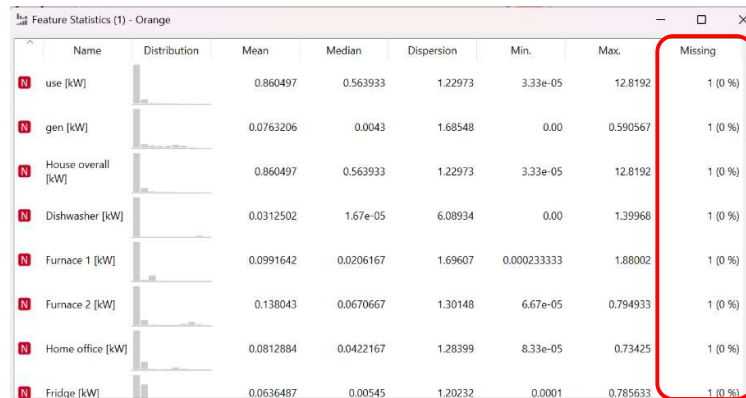


*Figure 2 Feature Statistics*

For the dataset provided for this course work, by understanding the dataset through feature statistics, It was clear that there were certain unknown values as shown in the fig 2.

With the help of the impute widget, we can clear out the missing values so that it may not affect the pipeline late on while the regression has to be performed.

Also, the meta-attribute Time is irrelevant to the energy generation. It is not necessary to edit the domain to convert this attribute into a feature.
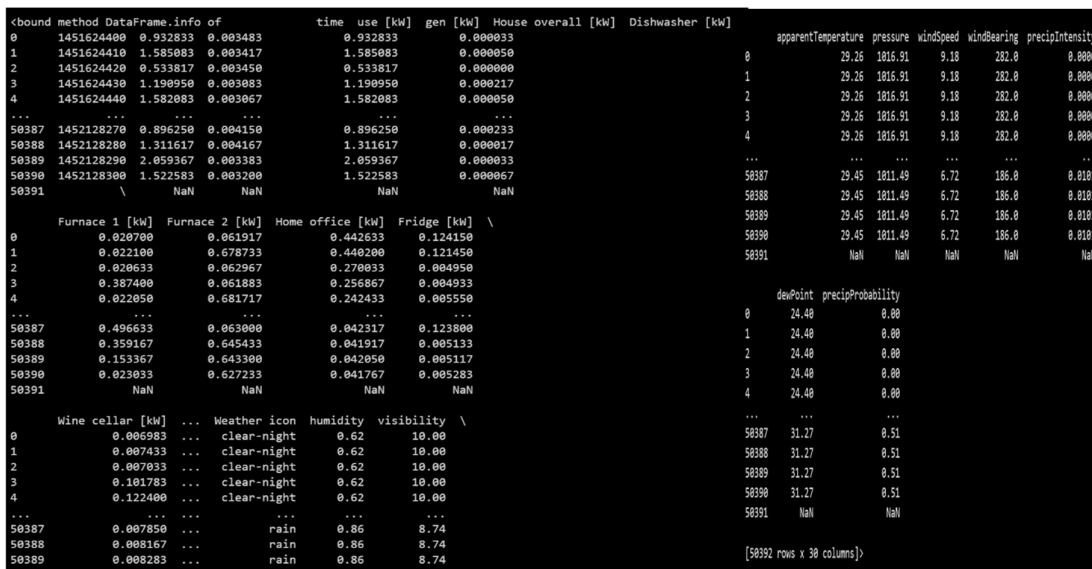
The current data description is given Figure 3.



*Figure 3 Data Description*

Once the data pre – processing is complete. The next step is to figure out the correlation between target feature and the rest of features.

## 4. Correlation

Machine learning models can be classified as good or based on the data that has been provided. The selection features contribute most to the quality of resulting model. The selection of such features that helps in predicting the variable more accurately is known as Feature selection.
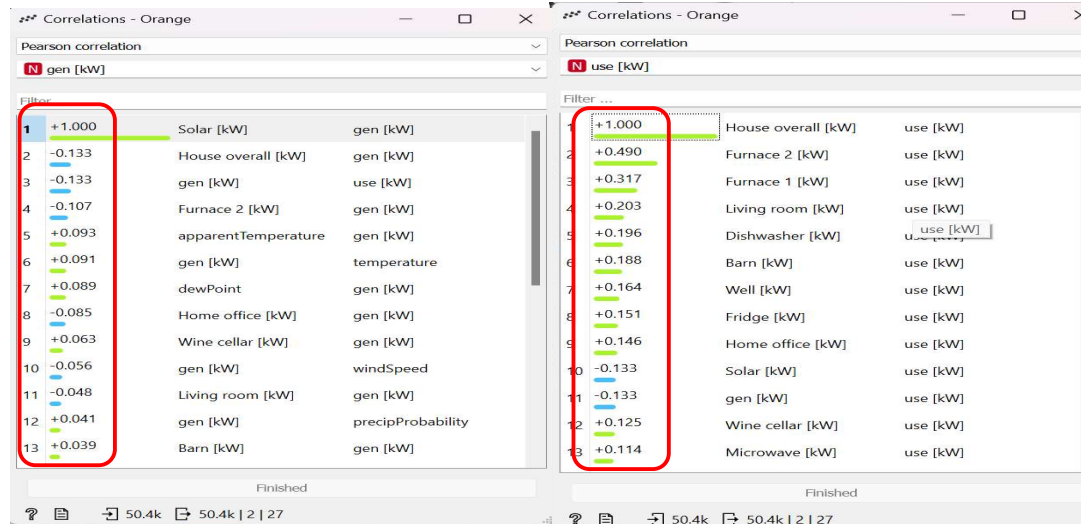


*Figure 4 Correlation*

Feature selection is crucial task which is based on configuring the relation between the available features. With help of such relationships, one can identify the features important based on the target model.

In case of Energy dataset, the correlations between features are given out in Fig (4). The marked values display the effect of the features on the target value that the use [kW]. But simply based on the correlation, it cannot be figured out, which features to drop and which to keep.
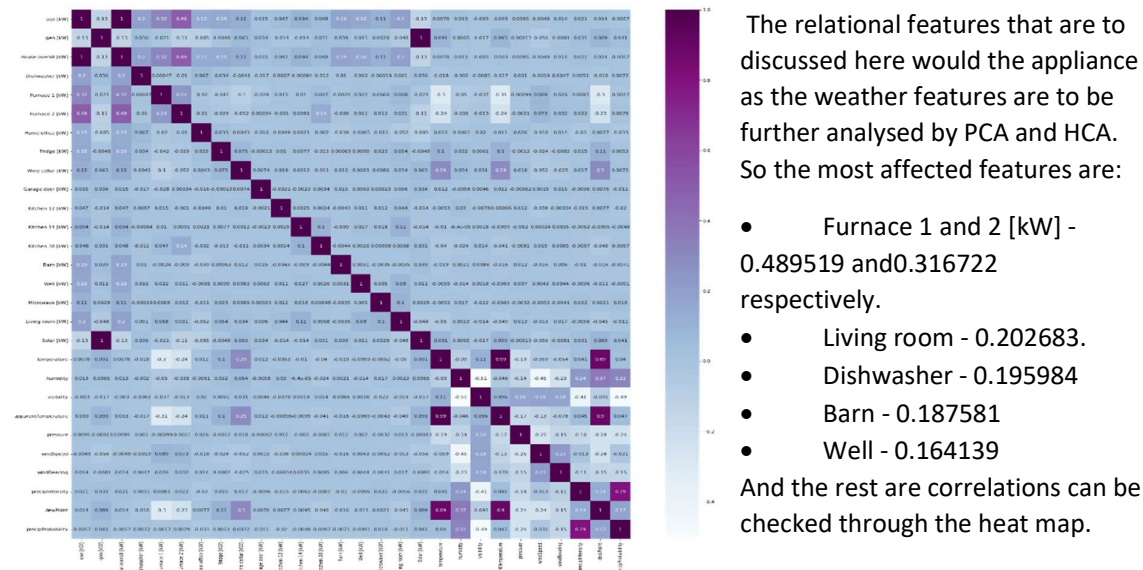


*Figure 5 Heat Map for Correlation*

The relational features that are to discussed here would the appliance as the weather features are to be further analysed by PCA and HCA. So the most affected features are:

- Furnace 1 and 2 [kW] - 0.489519 and 0.316722 respectively.
- Living room - 0.202683.
- Dishwasher - 0.195984
- Barn - 0.187581
- Well - 0.164139

And the rest are correlations can be checked through the heat map.

For the consideration of the weather features , Principal component analysis and hierarchical clustering analysis has to taken into account.  The pipeline for the has been given in the Figure 6.
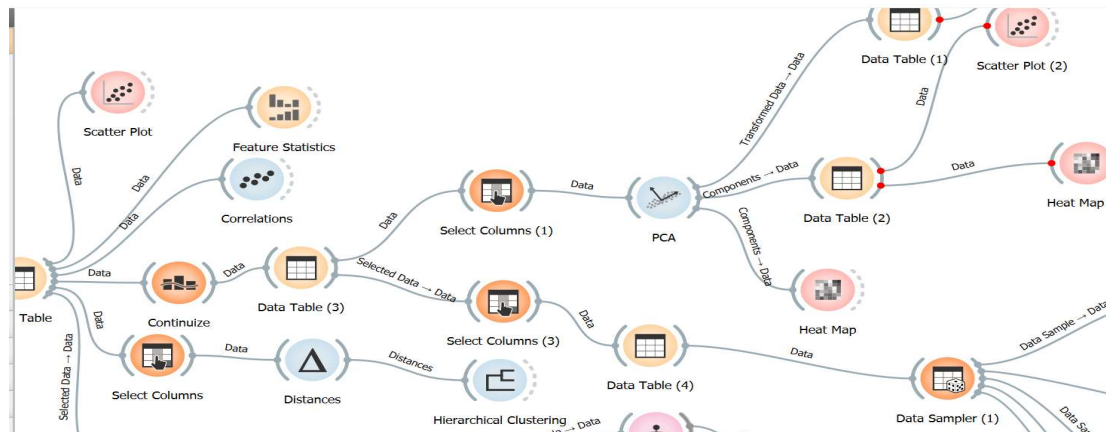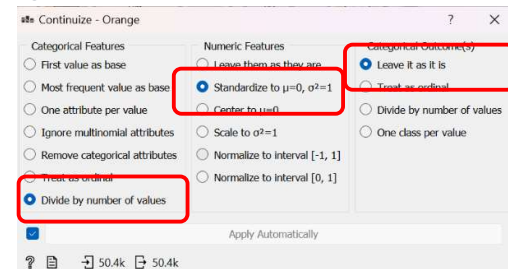


*Figure 6 Pipeline for Principal Component analysis and Hierarchical Clustering analysis*

From the pipeline, it is clearly understandable that to conduct to analysis, two branches are to be made. For the PCA, Continuize widget is specifically relevant as it determines the basic underlying conditions for the analysis. The condition for the same can be viewed in Figure 7.

*Figure 7 Continuize*



Once the underlying condition have been determined, the relevant column, i.e., the weather features have been selected with the help of the select column widget.  By the help of select column widget, only weather features have been selected to

confirm the variance amongst them. Next is to apply the PCA widget and check for the maximum variance with minimum number of components.

The variance can seen in the Figure 8.  It can be seen that with only 4 components the variance is 90% which is sufficient in terms of dimensionality reduction.  While reducing the PCs, the data presented with is shown in Figure 9.
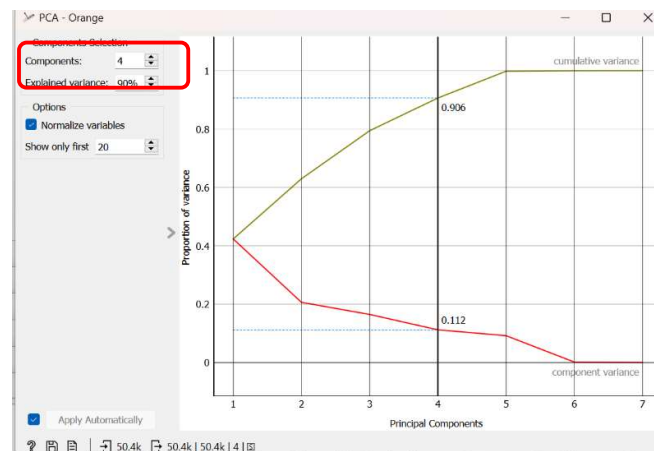


*Figure 8 Principal Component Analysis*



| | components | variance | pparentTemperatur | pressure | humidity | precipIntensity | dewPoint | temperature | windBearing |
|---|---|---|---|---|---|---|---|---|---|
| 1 | PC1 | 0.42327 | -0.5593 | 0.183764 | -0.101126 | -0.0975014 | -0.563509 | -0.556089 | 0.0830928 |
| 2 | PC2 | 0.206214 | -0.196018 | -0.164263 | 0.676281 | 0.516117 | 0.0929727 | -0.223158 | -0.390409 |
| 3 | PC3 | 0.164706 | -0.0745343 | -0.70251 | -0.0266611 | 0.2649 | -0.0526886 | -0.047983 | 0.651892 |
| 4 | PC4 | 0.112001 | 0.0779153 | 0.197879 | -0.541355 | 0.780707 | -0.165518 | 0.0978738 | -0.123407 |

Info
4 instances (no missing data)
7 features
No target variable.
2 meta attributes
Variables

*Figure 9 Analysed Data for PCA*

The above features have been reduced as the **Wind speed, Visibility and the precipitation probability** were inversely correlated to the to the target feature. Hence, they were kept out the analysis.

Also, the heat map of the data configured by the system should give out the additional data that is required to confirm the features to be added to make impact on the regression model of the system. This heat map is shown in Figure 9.
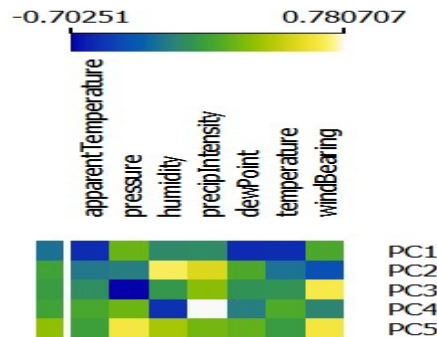


*Figure 10 Heat Map for PCA*

As for the Hierarchical Clustering, It is important to calculate the distances so as to confirm the relevance so that the features could be decided. Hence, the distances are calculated with the help of Distance widget wherein the metric unit is Euclidean. Once the distances are calculated, the clustering analysis is the next step. The output for the analysis is shown in Figure 10.
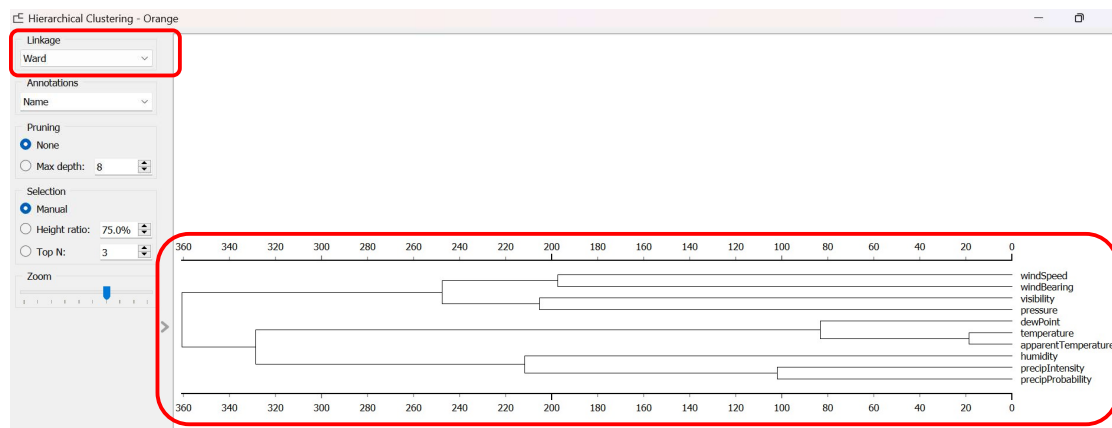


*Figure 10 Hierarchical Clustering Analysis*

With the help of the PCA and Hierarchical Clustering analysis, it is clear that the features with most weightage are given as:

- Pressure
- Wind Bearing
- Dewpoint
- Precipitation Intensity

Since, these features have been selected through above analysis, they will be included in the regression model further ahead.

## 5. Regression and Classification Model

➢ Regression model

Based on the domain analysis and the dataset's correlational study, it would be preferable to use multiple regression and cross validate them to obtain an accurate result for the system/pipeline. However, using the feature engineering done on the pipeline is not required.

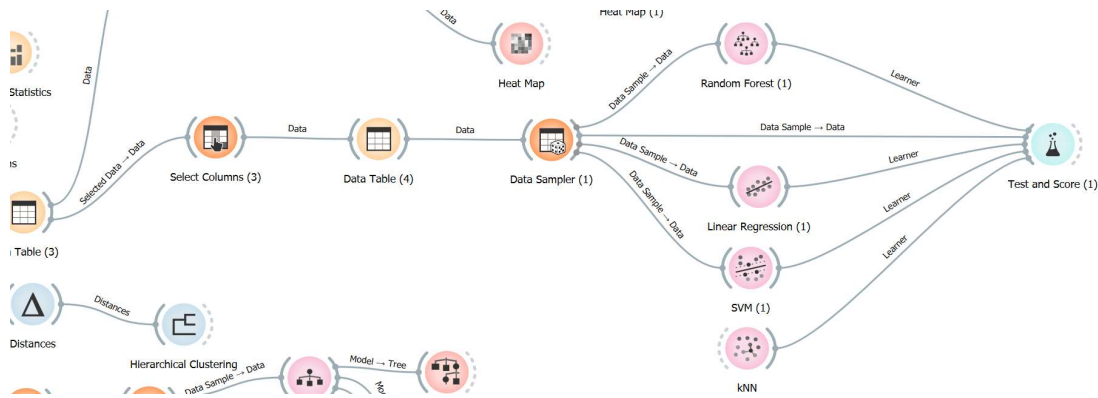For the regression model to be applicable, the pipeline would look like:



*Figure 11 Regression model pipeline*

As per the above pipeline, selecting the column that have an impact on the regression model based domain analysis and correlational study is the first step.

Through the PCA and HCA, the major impacting weather features have been added that can be viewed in the Figure 12.

The features that have been ignored are the result of the PCA , HCA and the correlational analsyis done earlier.

After that, all that remains is to divide the data into test and train datasets and to run the regression models through the Test and Score to determine which one is a better fit for the machine learning model.. The accuracy are visible in the Test and Score widget under the R2 column.
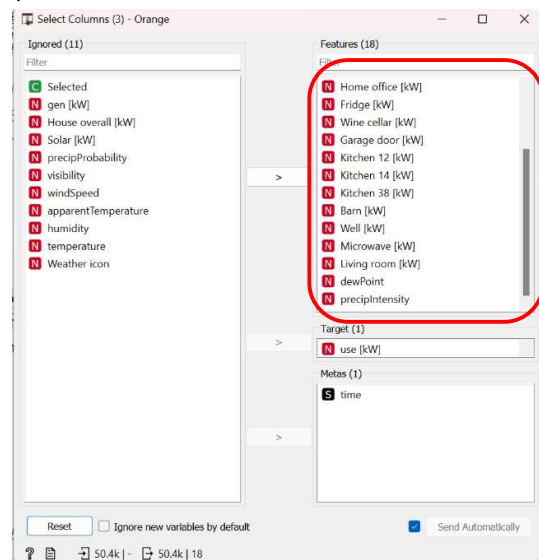


*Figure 12  Select column*

The Test and score stated that the random forest regression model is best fit based on the learning curve while linear regression and SVM seems to be a under fit. Although KNN seems to be a good fit as well but Random Forest regression show better accuracy. The values form orange tools are shown in Figure 13.

The above regression model results show that the Random Forest regression model is significantly more effective for the dataset based on Energy Consumption.

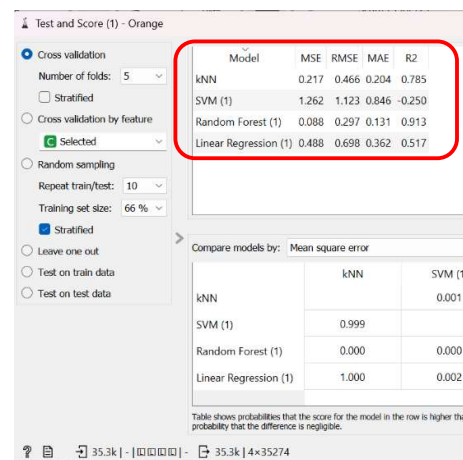Next Step would be to classify the system to figure out which data lie under High, Medium and Low usage.



| Model | MSE | RMSE | MAE | R2 |
|---|---|---|---|---|
| kNN | 0.217 | 0.466 | 0.204 | 0.785 |
| SVM (1) | 1.262 | 1.123 | 0.846 | -0.250 |
| Random Forest (1) | 0.088 | 0.297 | 0.131 | 0.913 |
| Linear Regression (1) | 0.488 | 0.698 | 0.362 | 0.517 |

*Figure 13 Regression Model result*

- Classification model

Classification models decide whether to categorise usage based on prior analysis performed through PCA, HCA, and domain analysis. However, in order to complete the classification model, feature engineering is required so that data can be studied easily.

In this case, the target feature must be converted to classification of 3: Low, Medium and High. Also, from the correlational analysis, it was observed that features Furnace 1 [kW] and Furnace 2 [kW] as well as Kitchen 12, Kitchen 14 and Kitchen 34 can be combine to a simplify the data classification. Hence, the new creations are done with the help of feature constructor widget.
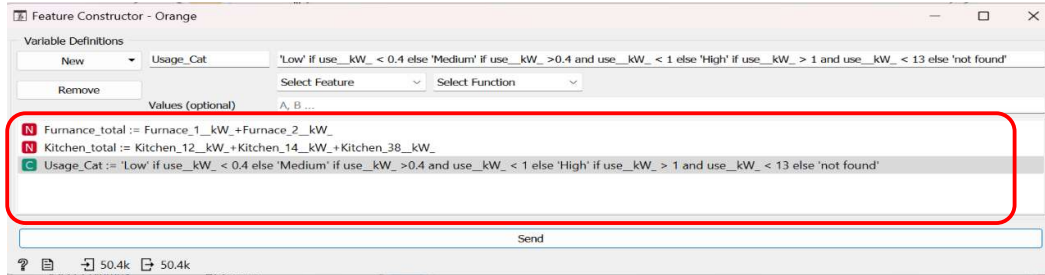


*Figure 14 Feature Construction*



*Figure 15 Data for Classification*

As presented in the Figure (14), new features have been created. All that is left to do is to use a sampler on the required data and take up the Decision Tree Classifier and Neural network.

Once the required columns have been chosen, sample and test data are generated to set up the system's classifier. In this case, a 70:30 ratio is used for the current pipeline.

When the pipeline's data sampler is prepared, the next step is to use this sample data for classifiers, with the remaining data being used to test the prediction. The classifier pipeline is as follows:

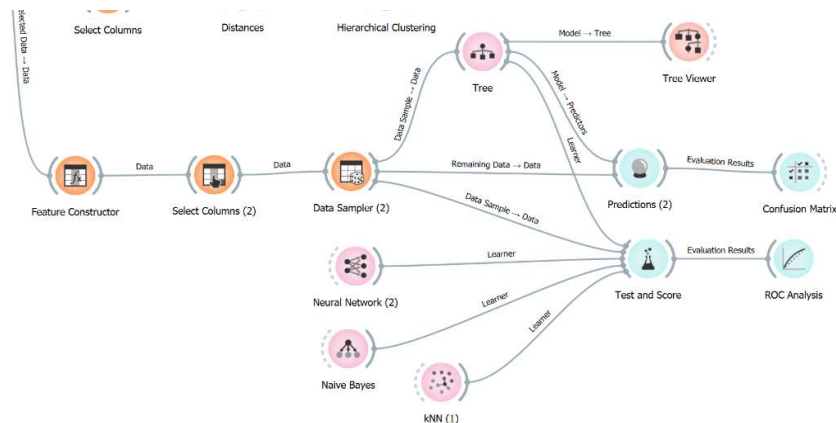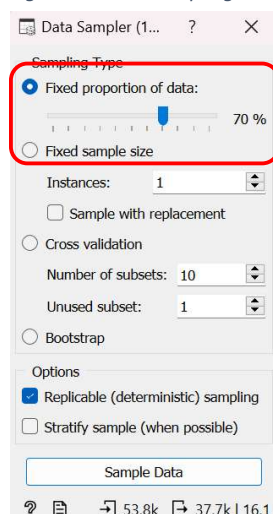*Figure 17 Data Sampling*





*Figure 16 Classification Model Pipeline*

From the above pipeline, the predictions are passed through the confusion matrix to make the sure of underfitting and over fitting for the classification.

The confusion matrix provides an accuracy to the predictions on the system so as to validate the pipeline.



*Figure 18 Predictions*

The Prediction shows the individual instance's predicted classification based on the pipeline into the three categories created under the feature Usage_Cat with presumptuous error.

Also, the precision for the classifier can also be seen in the same widget's window. As per the pipeline created, the system precision is given as 0.831

The prediction only classifies data, it is not capable of visualising it. As a result, a confusion matrix is used to provide the exact instances for the classification in terms of numbers and percentages.

The study of confusion matrix can be easily suggest out that accuracy of the pipeline and gives a better visualization of the dataset classification.



*Figure 19 Confusion Matrix*

A graphial representation for all the classifiers applied to the machine learning model to deduce which is more effective for the system. The ROC analysis shows the result in Figure 20.



*Figure 20 ROC Analysis`*

It is evident from the analysis that the Neural network is better option for the system as the output of the analysis shows that the curve for the neural network to be better.

- Model Complexity

Model complexity refers to the amount of detail or "wiggle room" that is built into a machine learning model. A model with high complexity will have many parameters or variables that can be adjusted, allowing it to capture a wide range of patterns and relationships in the data. For the dataset present and the models used the complexity given are below:

Neural networks, decision trees, and naive bayes are all machine learning algorithms that can be used for classification tasks. The complexity of each of these models can vary depending on a number of factors.

Neural networks are complex machine learning models that are composed of many interconnected nodes, or "neurons." The complexity of a neural network can be measured by the number of neurons and the number of connections, or weights, between them. As the number of neurons and connections increases, the complexity of the model also increases.

Decision trees are relatively simple models that consist of a series of branching decision rules As the number of decision rules and the depth of the tree increase, the complexity of the model also increases.

Naive bayes is a relatively simple probabilistic model that is based on the assumption of conditional independence among the features in the data. As the number of features and the number of possible values increases, the complexity of the model also increases.

K-nearest neighbors (kNN) is a relatively simple non-parametric model that is based on the concept of similarity between examples in the data. As the number of examples and features increases, the complexity of the model also increases.

## 6. Cross Validation

For a Machine learning model, a pipeline is created based on the requirement. But to actually verify the accuracy and the legitimacy of the model, cross validation is done. Essentially, it is method that evaluates and compares the learning algorithm by dividing data in two segments: one to train the model and the other one is used to test the predictive capabilities if the model such as accuracy, error, etc.

For this case, within the test and score result, it can be cross validated with help of K – fold cross validation where in the training data and testing data are split into multiple sets. Then for each set, the regression will take place and the average for all the regression will be shown as the final output.

The second method is by import the cross validation score from python. Which show the result for linear regression and random forest regression, SVM and kNN as:

As shown in Figure 21, cross vaildation for number of folds is given as 5 and the result shows that the random forest regressor is the best fit for the system.



*Figure 21 K fold Validation*

The validation prevents bias and is critical in maintaining the dataset required for the system to function.

It has an effect on the learning model because the test and train datasets must always be kept separate and have the exact ratio specified during time sampling.

## 7. Learning Curves – Underfitting and Overfitting

Learning curves are the plots that represent the performance and the impact of the model at the training data set increases. The two major causes for error in the pipeline presented are:

Bias: it describes the model so that it makes simpler assumptions such that, for the model, it is easier to approximate.

Variance: It describes the variability of the model's prediction with the change in dataset used for training the model.

What the learning curves helps in representing the exact regressor to be used for prediction as the system requires a good fit.

For this case, 3 regressors were used and trained on the same data sets. The regressors used were Linear regression, Random Forest regression and Neural Network. The plots are shown below:



Figure 22 Learning curve for Linear regression



Figure 23  Learning curve for Random Forest regression



Figure 24 Learning Curve for Neural network



Figure 25 Learning Curve for Decision tree

The reason for choosing the Random - forest regression as fit for the machine learning model can be easily deduced from the above plots for all the regression models.

For the learning curve of linear regression, the result output can be easily justified as the accuracy of the system not at par with the required model accuracy for machine learning model.

The neural network and random forest both have a good fit or the machine learning model, but the random forest seems to have the upper hand for cross validation curve. Hence, the better fit regressor for the system can be concluded to be the random forest.
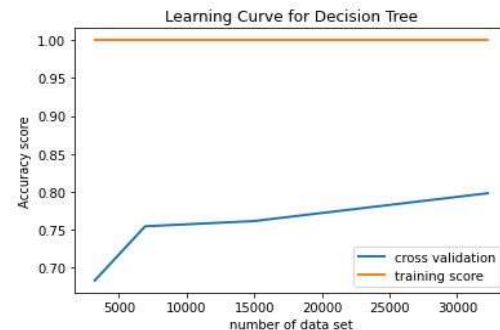
As per the data set used, once the classification is complete the dominance with highest used appliance in the three categories used can be studied with the help of the Linear projection widget. Applying the widget and using the dataset given that only includes the appliances. The projection is given in Figure 26.



*Figure 25 Linear projection*

The linear projection is a division of all the appliance based on the usage categories classified earlier by the feature engineering and along with that the study of the tree suggested that following appliance were being used during the prediction of the following categories:

- High

1. Barn [kW]

2. Well [kW]

- Medium

1. Living room [kW]

2. Home Office [kW]

- Low

1. Dishwasher [kW]

## 8. Mathematical Peculiarity for Machine learning Model

Mathematical Peculiarity for the system can be studied with the help of the output of the system and can be judged based on the accuracies as well the errors that the model present.

To study the peculiarity for the regression and classification models, we need a test and score with all the results of the process taken under single widget.

*Figure 26 Test and Score for Mathematical Peculiarity*

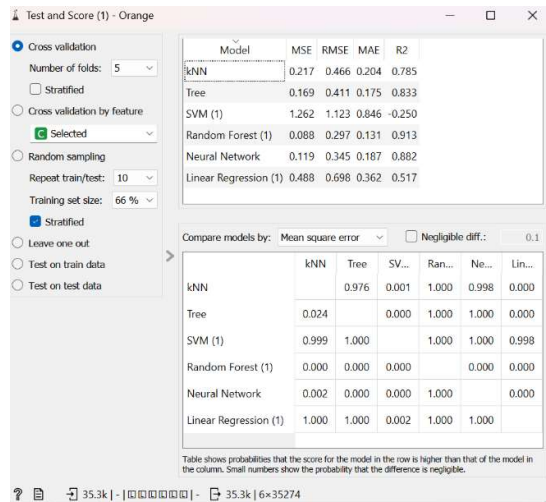It is clear that for linear regression the regression line always passes through the mean of the data. Also finding that the error present in the system is 0.698 which is too large for the model to be considered as a fit machine learning model.

As for the support vector machine regression, they use a kernel function to map the data from the original space into a higher-dimensional space. This allows the algorithm to find nonlinear decision boundaries that are more accurate than those found by other algorithms, such as logistic regression.

But the problem with SVM that the accuracy for the model is in negative which can never be used as a ML model.

One mathematical peculiarity of kNN is that it uses a distance metric to measure the similarity between data points. This distance metric can be Euclidean distance, Manhattan distance, or any other metric that quantifies the "distance" between two points in a given space.

For decision trees, the peculiarity is that they use a metric called entropy to measure the impurity or uncertainty of a node in the tree. Entropy is a measure of the amount of randomness or disorder in a system, and it is used to determine how well a node can be split into sub-nodes based on the values of the independent variables.

The mathematical peculiarity of random forests is that they use a technique called bootstrapping to train multiple decision trees on different subsets of the data. Bootstrapping is a sampling method that involves randomly selecting data points with replacement from the original dataset to create multiple smaller datasets.

- Evaluation matrix

For the given datasets, we have used 6 models in the pipeline. The resulting data given out is discussed in the report. Out of the six of them, but for the evaluation to be done. The matrix suggested that the SVM model could take place as the data set for Energy consumption and generation was too large and the SVM works perfectly for smaller data set.

Whereas the linear regression was not selected because of the accuracy it presented. It states that the prediction for the future use of the energy would only of 51% correct and that model, if used would have caused a lot of trouble for the pipeline.

The evaluation would have been between the rest of the models within the pipeline. kNN, Tree, Random Forest and neural network were taken into account as they were not underfit nor were they any close to being overfit. But the pipeline suggested that the better output was given by random forest of all the rest of them as it has the lowest RMSE and also the cross-validation curve is a rather better fit for random forest only.

The two models that were chosen out side the curriculum were kNN and Naive Bayes.

## 9. Conclusion

This course work consisted of creation of a complete pipeline for a data set based on study of 29 features of Energy consumption and generation. The process started with the domain analysis of the dataset. Once the domain analysis was complete. It gave enough information to construct a pipeline based on the features for prediction of prices of diamond and classifying them in three categories: low, Medium and High.

This process started with data cleaning and then moved forward to the system where it was required to understand which feature were to be kept and which were to be dropped. After the domain was complete. With the hep of the correlation study, PCA and HCA while keeping the target as use[kW], it was noticed that wind speed, visibility and precipitation probability had negative correlation. From the PCA and HCA, it was analysed which weather features to keep in the pipeline.

From here on, model was divided into two directions: First was creation a model that predict the usage of the energy. Second was the creation of a predictor that classified the usage of energy into three categories.

For the first part of the pipeline, a sampler was taken and the data was divided into a ratio of 70:30. The bigger set was used to train the model and the smaller set to test the regression models in place. It was through cross validation and learning curve found that the Random Forest regression model was a good fit for the machine learning model.

For the second part, feature engineer was done to categorise the class of the price and then through the sampler, a decision tree classification was done to achieve the confusion matrix that explain the prediction with error as well. Also, the pipeline showed a ROC analysis to confirm that the  Tree classifier was the correct choice.

With the completion of a machine learning model, the price was predicted and the classification was done which helped in creating an effecting ML Model.