

Tasks

The linear regression equation to model the Energy use:

$$\begin{aligned} \text{App} = & \beta_0 + \beta_{T1} * T1 + \beta_{RH_1} * RH_1 + \beta_{T2} * T2 + \beta_{RH_2} * RH_2 + \beta_{RH_5} * RH_5 + \\ & \beta_{T6} * T6 + \beta_{RH_6} * RH_6 + \beta_{T8} * T8 + \beta_{RH_8} * RH_8 + \beta_{T9} * T9 + \\ & \beta_{RH_9} * RH_9 + \beta_{Press} * Press + \beta_{RH_out} * RH_out + \beta_{Windspeed} * Windspeed \\ & + \beta_{Tdewpoint} * Tdewpoint \end{aligned}$$

The initial parameters are: Beta: 0.5 Number of iterations: 1000 Learning Rate: 0.01

Experiments

Experiment 1:

Linear Regression:

Experimenting with various paraments such as learning rate:

Choosing three learning Rates: [0.007, 0.005, 0.010]

Training Dataset:

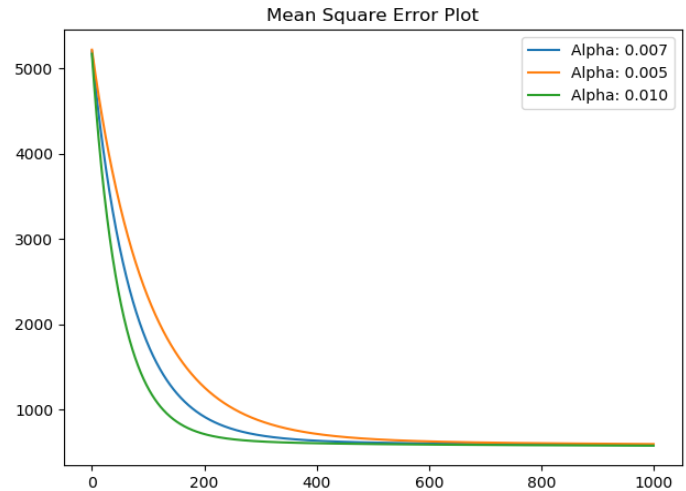
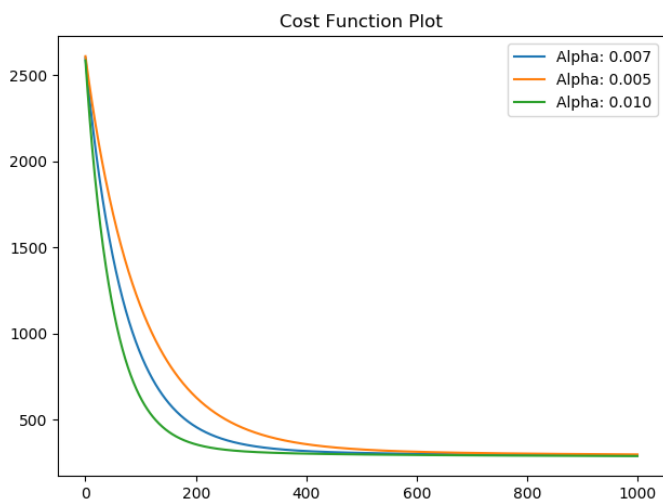
The Model was tested with varying learning rates and the best learning rate was found to be 0.01 as the number of iteration taken to complete the regression is the least while producing the best accuracy with its R square being the highest among the rest of the Learning Rates and the Cost Function being the Least.

```
Learning rate: 0.007, Iteration of Convergence: 630
Learning rate: 0.005, Iteration of Convergence: 882
Learning rate: 0.010, Iteration of Convergence: 440
```

Optimal Values for Training Data
Measures for Training Data

	CF	MSE	MAE	R ²
0.005	297.5127	595.0255	18.1099	0.2726
0.007	292.2345	584.4690	17.9992	0.2855
0.010	288.1791	576.3583	17.8707	0.2954

Varying Learning Rate



Test Dataset:

The varying learning rates are tested with the test data set, in this case too as you can see from the below images, 0.01 was the best Rsquare and the lowest Cost Function thus making it the best Rsquare

You can also observe that as the Learning Rate increases, the cost function decreases.

Measures for Test Data				
	CF	MSE	MAE	R ²
0.005	295.9177	591.8353	18.0782	0.2555
0.007	290.8510	581.7021	17.9728	0.2682
0.010	286.8035	573.6070	17.8356	0.2784

Logistic Regression:

Classes of the Dependent variable were classified using the median.

```
[0 if x <= 60 else 1
```

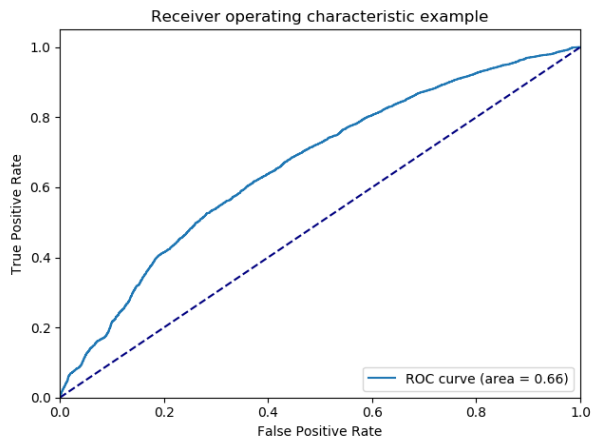
Logit Model applied on the Training Datasets has given me the following Metrics:

The learning rate was also varying with values such as 0.001, 0.005, 0.010, 0.050

The Value for the lowest Cost Function was 0.01 which converged after 2000 iterations.

Accuracy of the model in training was 26.59%

Area under the ROC curve: 0.663



Confusion Matrix

```
[[4832 2689]
```

```
 [1919 2877]]
```

Sensitivity: 0.5998748957464554

Specificity: 0.6424677569472145

Accuracy Score:

0.6258829260371844

Test:

The learning rate used when testing was 0.01, which produced the best cost function.

Accuracy of the Model on Test Data was 62.1%

Confusion matrix for Test

```
[[2034 1189]
```

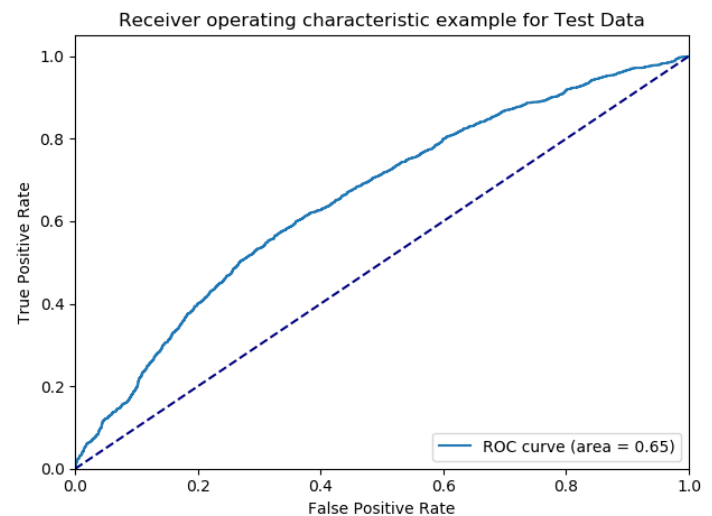
```
 [ 812 1245]]
```

Sensitivity: 0.6052503646086533

Specificity: 0.6310890474713

Accuracy Score:

0.6210227272727272



Experiment 2:

The different thresholds of cost function were implemented on the train dataset.

As the threshold decreases the accuracy of the model increases. The Cost Function decreases with decreasing threshold.

The iterations required to converge significantly decreased with Threshold value.

Train:

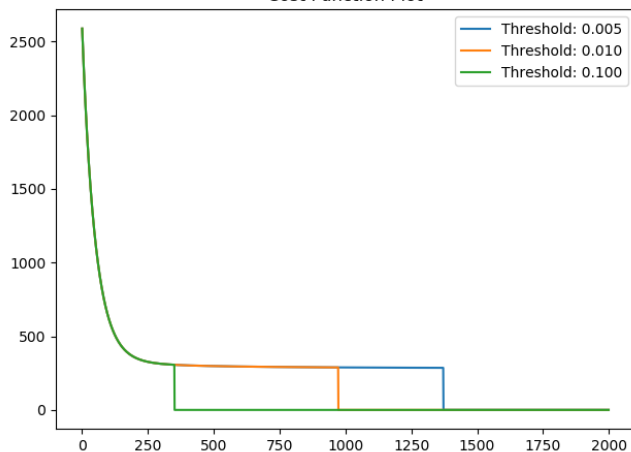
```
Learning rate: 0.010, Iteration of Convergence: 1371,  
Threshold: 0.005  
19:23:22.516451  
Learning rate: 0.010, Iteration of Convergence: 972,  
Threshold: 0.010  
19:23:40.444699  
Learning rate: 0.010, Iteration of Convergence: 350,  
Threshold: 0.100
```

Optimal Values for Training Data Experiment 2

Measures for Training Data Experiment 2

	CF	MSE	MAE	R ²
0.005	285.5945	571.1889	17.7862	0.3018
0.010	288.4433	576.8866	17.8793	0.2948
0.100	306.1593	612.3187	18.1112	0.2515

Cost Function Plot



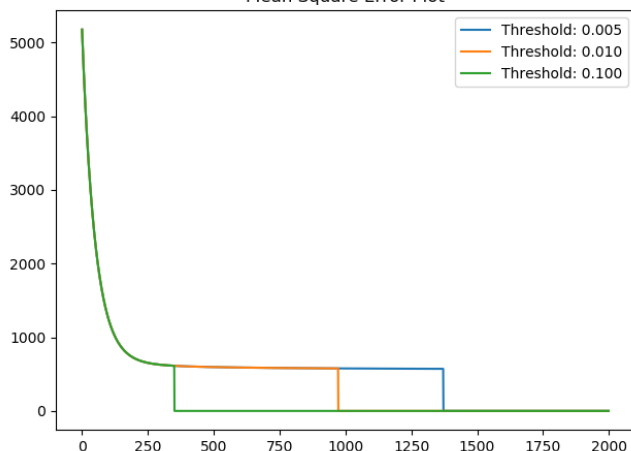
The graph shows the Cost Function with different thresholds

X- Axis: No. Of Iterations

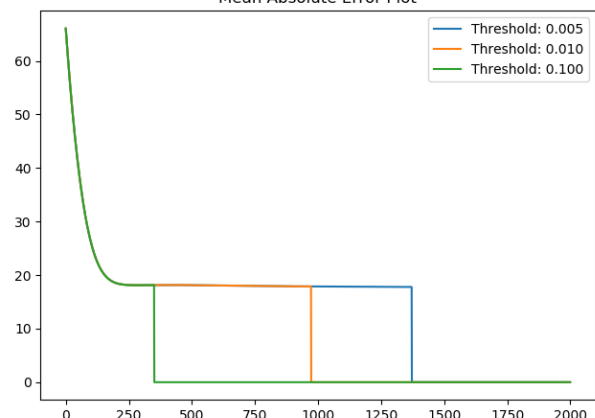
Y: Axis: Cost Function

As you can see in the graph the Cost Function speeds to an almost zero after it hits the Threshold.

Mean Square Error Plot



Mean Absolute Error Plot



The above graph also shows the best threshold and its graph plotted as a function of the number of iterations.

Test Data:

Using the different threshold values the measure of cost function and the r square value.

We can interpret that as we vary the threshold the cost function decreases with decrease in Threshold value

The R square on the other hand increases with decreasing threshold value, which holds true to the training data as well.

On varying different thresholds

Measures for Test Data

	CF	MSE	MAE	R ²
0.005	284.1312	568.2624	17.7356	0.2851
0.010	287.0727	574.1455	17.8455	0.2777
0.100	304.0246	608.0492	18.0629	0.2351

Experiment 3:

Linear regression:

Randomly picking the values:

Random Column picked are:

```
Index(['T2', 'T9', 'Visibility', 'T1', 'T4', 'RH_1', 'RH_7',  
'T7', 'T_out',  
      'T6'],  
      dtype='object')
```

The R square when randomly picked when compared to the 15 features picked shows a decrease i.e. the model is less accurate than the Optimized model.

Though the outliers were removed in both the cases, the features had a lot to do with the accuracy being down.

As the Visibility, T_out are some feature which has highly correlated and decrease the accuracy of the model.

Optimal Values for Training Data for Experiment 3

Measures for Training Data for Experiment 3

	CF	MSE	MAE	R ²
0.01	331.7969	663.5938	19.4359	0.1888

Testing on test data for Experiment 3

Measures for Test Data for experiment 3

	CF	MSE	MAE	R ²
0.01	329.883	659.7659	19.4493	0.17

Logit:

The Learning rate for the random variables was 0.01

The Rsquare like expected is less than the train dataset and the previously modeled data.

Optimal Values for Training Data for Experiment 3

Measures for Training Data for Experiment 3

	CF	MSE	MAE	R ²
0.01	331.7969	663.5938	19.4359	0.1888

Testing on test data for Experiment 3

Measures for Test Data for experiment 3

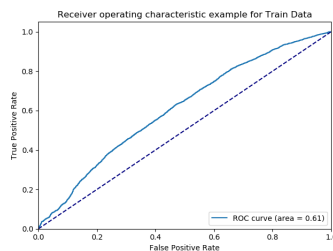
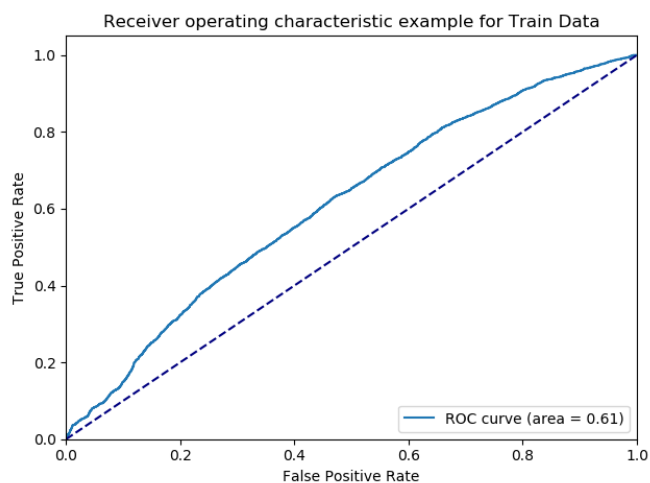
	CF	MSE	MAE	R ²
0.01	329.883	659.7659	19.4493	0.17

[4, 18, 24, 2, 8, 3, 15, 14, 20, 12]

The Random features picked are:

```
Index(['T2', 'T9', 'Visibility', 'T1', 'T4', 'RH_1', 'RH_7',  
      'T7', 'T_out',  
      'T6'],  
      dtype='object')
```

Logistic Regression



Experiment 4

Linear Regression:

Train':

Optimal Values for Training Data

Measures for Training Data

	CF	MSE	MAE	R ²
Page 9 of 10				

known

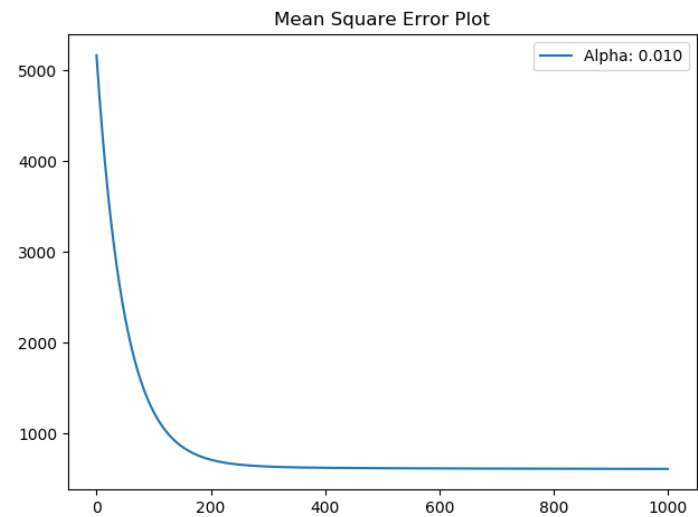
0.01	305.351	610.702	18.2914	0.2535
------	---------	---------	---------	--------

Test:

Testing on test data

Measures for Test Data

	CF	MSE	MAE	R ²
0.01	303.3763	606.7527	18.4842	0.2367



Mean Square Error Plot

